

Cheat Sheet

Contents

| | |
|---|----|
| Homework | 3 |
| 1: Speed up on multicore Processors | 3 |
| Assumptions | 3 |
| Task 1: Calculate performance on a dual-core chip | 3 |
| Task 2: Calculate performance on a quad-core chip | 4 |
| TA Response..... | 4 |
| 2: OpenMP Programming..... | 5 |
| 3: OpenMP Programming..... | 8 |
| Main.py..... | 8 |
| My Reduce function | 10 |
| My Reduce-Sequential function | 11 |
| 4: Google File System | 12 |
| 2: Cluster calculations..... | 13 |
| 5: MapReduce..... | 15 |
| OutDegree.java..... | 15 |
| WordCount.java..... | 16 |
| 6: Spark Programming | 17 |
| PageRank: | 17 |
| WordCount: | 18 |
| 7: Worse is better | 19 |
| What I learned..... | 19 |
| 8: Graph Partitioning | 20 |
| GraphPartitioner.py: | 20 |
| 9: Graph algorithms with GridGraph..... | 25 |
| Kcores.cpp: | 25 |
| PageRankDelta.cpp:..... | 26 |
| 10: Graph Mining..... | 28 |
| FSM.py | 28 |
| 11: Spark Streaming Top-K | 30 |

| | |
|-----------------------------|----|
| Top-k.py: | 30 |
| 12: Halide Scheduling | 32 |
| Dilated-conv.cpp: | 32 |

Homework

1: Speed up on multicore Processors

Assumptions

- 90% of a given program can be perfectly parallelized, the other 10% remains sequential
- We have a fixed power budget
- Total power is proportional to the square of frequency
- Performance is proportional to frequency
- The programs performance on a single-core chip is 1 solutions per seconds

Task 1: Calculate performance on a dual-core chip

On a single-core system, with the frequency f_1 , has the power usage of $P = cf_1^2$. Since the total power is fixed, the two dual-chip cores must share this power between them. We can therefore calculate the frequency f_2 for each of the cores:

$$\frac{P}{2} = cf_2^2$$

$$\frac{cf_1^2}{2} = cf_2^2$$

$$\frac{f_1}{\sqrt{2}} = f_2$$

Assuming that on the sequential sections of the program, only one core is working, we can therefore calculate the performance for the throughput for the parallelizable and sequential sections of the program

$$t_p = 2 * f_2 = 2 * \frac{f_1}{\sqrt{2}} = \sqrt{2} * f_1 = \sqrt{2} \text{ solutions/s}$$

$$t_s = 1 * f_2 = \frac{f_1}{\sqrt{2}} = \frac{1}{\sqrt{2}} \text{ solutions/s}$$

We can now calculate the total throughput, and therefore performance, of the program on the dual-core chip

$$t_{total} = 0.9 * t_p + 0.1 * t_s$$

$$t_{total} = 0.9 * \sqrt{2} \text{ solutions/s} + 0.1 * \frac{1}{\sqrt{2}} \text{ solutions/s}$$

$$t_{total} = \frac{19\sqrt{2}}{20} \text{ solutions/s}$$

$$t_{total} \approx 1,344 \text{ solutions/s}$$

Task 2: Calculate performance on a quad-core chip

As before, we can calculate the frequency f_4 for each of the four cores:

$$\frac{P}{4} = cf_4^2$$

$$\frac{cf_1^2}{4} = cf_4^2$$

$$\frac{f_1}{\sqrt{4}} = f_4$$

$$\frac{f_1}{2} = f_4$$

Again, we calculate the performance for the throughput for the parallelizable and sequential sections of the program

$$t_p = 4 * f_4 = 4 * \frac{f_1}{2} = 2f_1 = 2 \text{ solutions/s}$$

$$t_s = 1 * f_4 = \frac{f_1}{2} = \frac{1}{2} \text{ solutions/s}$$

We now calculate the total throughput

$$t_{total} = 0.9 * t_p + 0.1 * t_s$$

$$t_{total} = 0.9 * 2 \text{ solutions/s} + 0.1 * \frac{1}{2} \text{ solutions/s}$$

$$t_{total} = 1.85 \text{ solutions/s}$$

TA Response

Wrong math on speed averaging. Overall performance does not equal to the sum of the performance times the proportion of each part.

2: OpenMP Programming

```
#include <getopt.h>
#include <omp.h>
#include <stdio.h>
#include <stdlib.h>
#include <string>

#include <iostream>
#include <sstream>
#include <vector>

#include "./common/CycleTimer.h"
#include "./common/grade.h"
#include "./common/graph.h"
#include "page_rank.h"

#define USE_BINARY_GRAPH 1
#define PageRankDampening 0.3f
#define PageRankConvergence 1e-7d

// used for check correctness
void reference_serial_pageRank(Graph g, double *solution, double damping,
                                double convergence) {
    int numNodes = num_nodes(g);
    double equal_prob = 1.0 / numNodes;
    double *solution_new = new double[numNodes];
    double *score_old = solution;
    double *score_new = solution_new;
    bool converged = false;
    double broadcastScore = 0.0;
    double globalDiff = 0.0;
    int iter = 0;

    for (int i = 0; i < numNodes; ++i) {
        solution[i] = equal_prob;
    }
    while (!converged && iter < MAXITER) {
        iter++;
        broadcastScore = 0.0;
        globalDiff = 0.0;
        for (int i = 0; i < numNodes; ++i) {
            score_new[i] = 0.0;

            if (outgoing_size(g, i) == 0) {
                broadcastScore += score_old[i];
            }
            const Vertex *in_begin = incoming_begin(g, i);
            const Vertex *in_end = incoming_end(g, i);
            for (const Vertex *v = in_begin; v < in_end; ++v) {
                score_new[i] += score_old[*v] / outgoing_size(g, *v);
            }
            score_new[i] =
                damping * score_new[i] + (1.0 - damping) * equal_prob;
        }
        for (int i = 0; i < numNodes; ++i) {
            score_new[i] += damping * broadcastScore * equal_prob;
            globalDiff += std::abs(score_new[i] - score_old[i]);
        }
        converged = (globalDiff < convergence);
        std::swap(score_new, score_old);
    }
    if (score_new != solution) {
        memcpy(solution, score_new, sizeof(double) * numNodes);
    }
    delete[] solution_new;
}
```

```

int main(int argc, char **argv) {

    int num_threads = -1;
    std::string graph_filename;

    if (argc < 3) {
        std::cerr << "Usage: <path/to/graph/file> <manual_set_thread_count>\n";
        exit(1);
    }

    int thread_count = -1;
    if (argc == 3) {
        thread_count = atoi(argv[2]);
    }
    if (thread_count <= 0) {
        std::cerr << "<manual_set_thread_count> must > 0\n";
        exit(1);
    }
    graph_filename = argv[1];

    Graph g;

    printf("-----\n");
    printf("Running with %d threads\n", thread_count);
    printf("-----\n");
    printf("Loading graph...\n");

    if (USE_BINARY_GRAPH) {
        g = load_graph_binary(graph_filename.c_str());
    } else {
        g = load_graph(argv[1]);
        printf("storing binary form of graph!\n");
        store_graph_binary(graph_filename.append(".bin").c_str(), g);
        delete g;
        exit(1);
    }
    printf("\n");
    printf("Graph stats:\n");
    printf("  Filename: %s\n", argv[1]);
    printf("  Edges: %d\n", g->num_edges);
    printf("  Nodes: %d\n", g->num_nodes);

    bool pr_check = true;
    double *sol1;
    sol1 = (double *)malloc(sizeof(double) * g->num_nodes);
    double *sol2;
    sol2 = (double *)malloc(sizeof(double) * g->num_nodes);

    double pagerank_base;
    double pagerank_time;

    double ref_pagerank_base;
    double ref_pagerank_time;

    double start;
    std::stringstream timing;
    std::stringstream ref_timing;

    timing << "Threads Page Rank\n";
    ref_timing << "Serial Reference Page Rank\n";

    // Set thread count
    omp_set_num_threads(thread_count);

    // Run implementations
    start = CycleTimer::currentSeconds();
    pageRank(g, sol1, PageRankDampening, PageRankConvergence);
    pagerank_time = CycleTimer::currentSeconds() - start;
}

```

```

// Run reference implementation
start = CycleTimer::currentSeconds();
reference_serial_pageRank(g, sol2, PageRankDampening, PageRankConvergence);
ref_pagerank_time = CycleTimer::currentSeconds() - start;

printf("-----\n");
std::cout << "Testing Correctness of Page Rank\n";
if (!compareApprox(g, sol2, sol1)) {
    pr_check = false;
}

if (!pr_check)
    std::cout << "Your Page Rank is not Correct" << std::endl;
else
    std::cout << "Your Page Rank is Correct" << std::endl;

char buf[1024];
char ref_buf[1024];
sprintf(buf, "%4d: %.6f s\n", thread_count, pagerank_time);
sprintf(ref_buf, " 1: %.6f s\n", ref_pagerank_time);

timing << buf;
ref_timing << ref_buf;

printf("-----\n");
std::cout << "Serial Reference Summary" << std::endl;
std::cout << ref_timing.str();

printf("-----\n");
std::cout << "Timing Summary" << std::endl;
std::cout << timing.str();

printf("-----\n");

delete g;
return 0;
}

```

3: OpenMP Programming

Main.py

```
#include "your_reduce.h"
#include <cassert>
#include <cstdio>
#include <cstdlib>
#include <cstring>
#include <ctime>
#include <random>
#include <sys/time.h>

#define MAX_LEN 268435456

// return the time in the unit of us
static long get_time_us() {
    struct timeval my_time;
    gettimeofday(&my_time, NULL);
    long runtime_us = 1000000 * my_time.tv_sec + my_time.tv_usec;
    return runtime_us;
}

int main(int argc, char *argv[]) {
    int size, rank, provided;

    MPI_Init_thread(&argc, &argv, MPI_THREAD_MULTIPLE,
                    &provided); // enable multi-thread support (for Bonus)
    assert(provided == MPI_THREAD_MULTIPLE);

    MPI_Comm_size(MPI_COMM_WORLD, &size);
    MPI_Comm_rank(MPI_COMM_WORLD, &rank);

    int *a, *b; // the array used to do the reduction

    int *res; // the array to record the result of YOUR_Reduce
    int *res2; // the array to record the result of MPI_Reduce

    long count;
    long begin_time, end_time, use_time,
        use_time2; // use_time for YOUR_Reduce & use_time2 for MPI_Reduce

    int i;

    // initialize
    a = (int *)malloc(MAX_LEN * sizeof(int));
    b = (int *)malloc(MAX_LEN * sizeof(int));
    res = (int *)malloc(MAX_LEN * sizeof(int));
    res2 = (int *)malloc(MAX_LEN * sizeof(int));
    memset(a, 0, MAX_LEN * sizeof(int));
    memset(b, 0, MAX_LEN * sizeof(int));
    memset(res, 0, MAX_LEN * sizeof(int));
    memset(res2, 0, MAX_LEN * sizeof(int));

    std::mt19937 rng;
    rng.seed(time(NULL)); // seed to generate the array randomly

    for (count = 1; count <= MAX_LEN;
         count *= 16) // length of array : [ 1 16 256 4'096 65'536 1'048'576
                  // 16'777'216 268'435'456 ]
    // do not report results for length 1
    {
        // the element of array is generated randomly
        for (i = 0; i < count; i++) {
            b[i] = a[i] = rng() % MAX_LEN;
        }
    }
}
```

```

// MPI_Reduce and then print the usetime, the result will be put in
// res2[]
MPI_Barrier(MPI_COMM_WORLD);
begin_time = get_time_us();
MPI_Reduce(a, res2, count, MPI_INT, MPI_SUM, 0, MPI_COMM_WORLD);
MPI_Barrier(MPI_COMM_WORLD);
end_time = get_time_us();
use_time2 = end_time - begin_time;
if (rank == 0)
    printf("%ld int use_time : %ld us [MPI_Reduce]\n", count,
           use_time2),
    fflush(stdout);

// YOUR_Reduce and then print the usetime, the result should be put in
// res[]
MPI_Barrier(MPI_COMM_WORLD);
begin_time = get_time_us();

YOUR_Reduce(b, res, count);

MPI_Barrier(MPI_COMM_WORLD);
end_time = get_time_us();
use_time = end_time - begin_time;
if (rank == 0)
    printf("%ld int use_time : %ld us [YOUR_Reduce]\n", count,
           use_time),
    fflush(stdout);

// check the result of MPI_Reduce and YOUR_Reduce
if (rank == 0) {
    int correctness = 1;
    for (i = 0; i < count; i++) {
        if (res2[i] != res[i]) {
            correctness = 0;
        }
    }
    if (correctness == 0)
        printf("WRONG !!!\n"), fflush(stdout);
    else
        printf("CORRECT !\n"), fflush(stdout);
}
}

MPI_Finalize();

return 0;
}

```

My Reduce function

```
#include <mpi.h>
#include <cstring>
#include <stdio.h>

void YOUR_Reduce(const int *sendbuf, int *recvbuf, int count) {
    int rank, size;

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    // Initialize recvbuf with the values from sendbuf
    memcpy(recvbuf, sendbuf, count * sizeof(int));

    // Allocate temp_buffer once
    int* temp_buffer = new int[count];

    // Binary tree reduction
    for (int step = 1; step < size; step *= 2) {
        if (rank % (2 * step) == 0) {
            // Root process collects data
            if (rank + step < size) {
                MPI_Recv(temp_buffer, count, MPI_INT, rank + step, 0,
MPI_COMM_WORLD, MPI_STATUS_IGNORE);

                // Combine results using a loop
                for (int i = 0; i < count; i++) {
                    recvbuf[i] += temp_buffer[i];
                }
            }
        } else if (rank % step == 0) {
            // Send to parent
            MPI_Send(recvbuf, count, MPI_INT, rank - step, 0, MPI_COMM_WORLD);
            break;
        }
    }

    // Clean up
    delete[] temp_buffer;
}
```

My Reduce-Sequential function

```
#include <mpi.h>
#include <cstring>
#include <stdio.h>
#include <omp.h> // Include OpenMP header

void YOUR_Reduce(const int *sendbuf, int *recvbuf, int count) {
    int rank, size;

    MPI_Comm_rank(MPI_COMM_WORLD, &rank);
    MPI_Comm_size(MPI_COMM_WORLD, &size);

    // Initialize recvbuf with the values from sendbuf
    memcpy(recvbuf, sendbuf, count * sizeof(int));

    // Allocate temp_buffer once
    int* temp_buffer = new int[count];

    // Binary tree reduction
    for (int step = 1; step < size; step *= 2) {
        if (rank % (2 * step) == 0) {
            // Root process collects data
            if (rank + step < size) {
                MPI_Request request;

                // Start non-blocking receive
                MPI_Irecv(temp_buffer, count, MPI_INT, rank + step, 0,
                          MPI_COMM_WORLD, &request);

                // Wait for the receive to complete
                MPI_Wait(&request, MPI_STATUS_IGNORE);

                // Combine results in parallel
                #pragma omp parallel for
                for (int i = 0; i < count; i++) {
                    recvbuf[i] += temp_buffer[i];
                }
            }
        } else if (rank % step == 0) {
            // Send to parent using non-blocking send
            MPI_Request request;
            MPI_Isend(recvbuf, count, MPI_INT, rank - step, 0, MPI_COMM_WORLD,
                      &request);
            break;
        }
    }

    // Clean up
    delete[] temp_buffer;
}
```

4: Google File System

1: GFS Questions

1.1: How does the master node get the locations of each chunks at startup?

At startup, the master node does not store the chunk location information persistently. Instead, it retrieves the location of each chunk by polling all chunkservers. Each chunkserver reports the chunks it holds, and the master node updates its information accordingly. Additionally, whenever a chunkserver joins the cluster, the master node updates the chunk locations. The master keeps this information updated by sending periodic HeartBeat messages, making sure it is updated on chunkplacement and the chunkserver status. This approach ensures that the master always has up-to-date information about chunk locations in the system.

1.2: What is the benefit of this approach comparing with the approach that the master persists this information?

The main benefit of this approach is the reduction in complexity related to maintaining consistency between the master and the chunkservers. Persisting chunk location information would require the system to handle various events such as chunkserver failures, renaming, and rejoining, which can lead to stale or inconsistent information. By polling the chunkservers at startup and using regular HeartBeat messages, the master node avoids the issues of synchronization and ensures accurate, up-to-date chunk location information at all times. This also simplifies the system's design, making it more robust against common failures in a distributed environment.

2: Cluster calculations

We assume a cluster of 1000 servers, each server having 10 disks with 10TB storage capacity and 100MB/s I/O bandwidth per disk. The servers are connected by a 1Gbps (125Mbps) ethernet cables, as nothing else is written I assume this bandwidth is per machine and not the total transfer cap in the system.

2.1: What is the minimum time required to recovery a node failure

The total I/O bandwidth of all disks in a node is:

$$\text{Total bandwidth} = 10\text{disks} \times 100\text{MB/s} = 1000\text{MB/s}$$

The network bandwidth available per node is:

$$\text{Network bandwidth} = 1\text{Gbps} = 1024\text{Mbps} = 128\text{MB/s}$$

Since the network bandwidth is lower than the total I/O bandwidth of the disks, the recovery time will be limited by the network speed. Furthermore, we assume that since we are calculating the minimum time required, that all other 999 nodes will use their resources to assist upon a node failure. Assuming that the 999 remaining servers work in perfect parallel to recover the data from the failed node, the total network bandwidth available is

$$\text{Total bandwidth} = 999 \times 128\text{MB/s} = 127872\text{MB/s}$$

To transfer 100 TB of data at this rate, the time required is:

$$\text{Time to recover} = \frac{100 \times 1024 \times 1024\text{MB}}{127872\text{MBs}} = 820\text{s} = 13,667\text{m}$$

Thus, the minimum time required to recover a node failure, assuming all other servers participate in the recovery, is approximately 14 minutes.

2.2: Time to recover a failure node with throttled recovery bandwidth

Since the recovery traffic is throttled to 100 Mbps per machine, roughly a tenth of the original 1Gbps, we would expect a recovery time ten times as long, as the recovery time is directly proportional to the network. Here we still assume that all other 999 nodes, as nothing else is stated in the assignment. We also assume all other requirements are similar. Below are my calculations:

$$\text{Total bandwidth} = 999 \times \frac{100\text{ Mbps}}{8\text{ MB/Mb}} = 999 \times 12.5\text{ MB/s} = 12487.5\text{ MB/s}$$

At this throttled rate, the time required to recover 100 TB of data is:

$$\text{Time to recover} = \frac{100 \times 1024 \times 1024\text{ MB}}{12487.5\text{ MB/s}} = 8397\text{ seconds}$$

Converting this to hours:

$$\text{Time to recover} = \frac{8392\text{ seconds}}{3600\text{ seconds/h}} \approx 2.33\text{ hours.}$$

Thus, the time required to recover a node failure when the recovery traffic is throttled is approximately 2 hours and 20 minutes.

Q3: How many server failures are likely to occur in a year in this cluster? What is the mean time between node failures in this cluster?

We know that each server node has a mean time between failures (MTBF) of 10,000 hours.

$$\frac{24 \text{ h/day} \times 365 \text{ days/year}}{10,000 \text{ h MTBF}} = \frac{8,760 \text{ h/year}}{10,000 \text{ h/failure}} = 0.876 \text{ failures/year/server.}$$

As this failure rate is independent for each server node, we have to multiply this by the total amount of servers in the cluster to estimate the expected number of failures per year in the cluster:

$$1000 \text{ servers} \times 0.876 \text{ failures/year/server} = 876 \text{ failures/year.}$$

The mean time between failures (MTBF) for the entire cluster is therefore:

$$\frac{8,760 \text{ hours/year}}{876 \text{ failures/year}} = 10 \text{ hours.}$$

Thus, the cluster is expected to experience one server failure approximately every 10 hours.

Q4: What is the implication of the number of replicas used in GFS based on the results from Q2 and Q3?

The comparison between the recovery time (Q2) and the mean time between node failures (Q3) highlights the critical importance of replication in GFS. With recovery taking approximately 2.33 hours when throttled (as calculated in Q2), and with the cluster experiencing a node failure approximately every 10 hours (from Q3), it's evident that replication is essential to prevent data loss. The default number of replicas in GFS is three, which ensures redundancy and availability of data even when nodes fail.

The chances of multiple replicas experiencing data loss due to node failures are slim, and the chances of this is greatly reduced for each copy. Still, increasing the number leads to a reduction in performance as more updates have to be completed for each chunk write. Still, the default of three replicas strikes a good balance for most use cases. It provides sufficient fault tolerance while keeping the storage overhead manageable. Increasing the number of replicas could provide additional safety in environments with higher failure rates, but this would come at the cost of additional storage requirements. Conversely, reducing the number of replicas would expose the system to an increased risk of data loss in the event of multiple simultaneous failures, which GFS is designed to avoid.

5: MapReduce

OutDegree.java

```
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;
import org.apache.hadoop.util.GenericOptionsParser;

public class OutDegree {

    public static class OutDegreeMapper
        extends Mapper<Object, Text, Text, IntWritable> {

        private final static IntWritable one = new IntWritable(1);
        private Text node = new Text();

        public void map(Object key, Text value, Context context
                        ) throws IOException, InterruptedException {
            StringTokenizer itr = new StringTokenizer(value.toString());
            itr.nextToken(); // Skip the first token ("a")
            if (itr.hasMoreTokens()) {
                String sourceNode = itr.nextToken(); // The source node (e.g., "a")
                node.set(sourceNode);
                // Emit the source node with a count of 1 for each outgoing edge
                context.write(node, one);
            }
        }
    }

    public static class OutDegreeReducer
        extends Reducer<Text, IntWritable, Text, IntWritable> {
        private IntWritable result = new IntWritable();

        public void reduce(Text key, Iterable<IntWritable> values,
                           Context context
                           ) throws IOException, InterruptedException {
            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get(); // Sum up all counts for each node
            }
            result.set(sum);
            context.write(key, result); // Emit the node and its total out-degree
        }
    }

    public static void main(String[] args) throws Exception {
        Configuration conf = new Configuration();
        String[] otherArgs = new GenericOptionsParser(conf, args).getRemainingArgs();
        if (otherArgs.length < 2) {
            System.err.println("Usage: outdegree <in> <out>");
            System.exit(2);
        }
        Job job = new Job(conf, "outdegree");
        job.setJarByClass(OutDegree.class);
        job.setMapperClass(OutDegreeMapper.class);
        job.setCombinerClass(OutDegreeReducer.class);
        job.setReducerClass(OutDegreeReducer.class);
        job.setNumReduceTasks(1);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        for (int i = 0; i < otherArgs.length - 1; ++i) {
            FileInputFormat.addInputPath(job, new Path(otherArgs[i]));
        }
        FileOutputFormat.setOutputPath(job, new Path(otherArgs[otherArgs.length - 1]));
        System.exit(job.waitForCompletion(true) ? 0 : 1);
    }
}
```

WordCount.java

```
import java.io.IOException;
import java.util.StringTokenizer;

import org.apache.hadoop.conf.Configuration;
import org.apache.hadoop.fs.Path;
import org.apache.hadoop.io.IntWritable;
import org.apache.hadoop.io.Text;
import org.apache.hadoop.mapreduce.Job;
import org.apache.hadoop.mapreduce.Mapper;
import org.apache.hadoop.mapreduce.Reducer;
import org.apache.hadoop.mapreduce.lib.input.FileInputFormat;
import org.apache.hadoop.mapreduce.lib.output.FileOutputFormat;

public class WordCount {

    public static class TokenizerMapper extends Mapper<Object, Text, Text, IntWritable> {

        private final static IntWritable one = new IntWritable(1);
        private Text word = new Text();

        public void map(Object key, Text value, Context context) throws IOException, InterruptedException {
            StringTokenizer itr = new StringTokenizer(value.toString());
            while (itr.hasMoreTokens()) {
                word.set(itr.nextToken());
                context.write(word, one);
            }
        }
    }

    public static class IntSumReducer extends Reducer<Text, IntWritable, Text, IntWritable> {
        private IntWritable result = new IntWritable();

        public void reduce(Text key, Iterable<IntWritable> values, Context context)
                throws IOException, InterruptedException {
            int sum = 0;
            for (IntWritable val : values) {
                sum += val.get();
            }
            result.set(sum);
            context.write(key, result);
        }
    }

    public static int run(String input, String output) throws Exception {
        Configuration conf = new Configuration();
        Job job = Job.getInstance(conf, "word count");
        job.setJarByClass(WordCount.class);
        job.setMapperClass(TokenizerMapper.class);
        job.setCombinerClass(IntSumReducer.class);
        job.setReducerClass(IntSumReducer.class);
        job.setOutputKeyClass(Text.class);
        job.setOutputValueClass(IntWritable.class);
        FileInputFormat.addInputPath(job, new Path(input));
        FileOutputFormat.setOutputPath(job, new Path(output));
        return job.waitForCompletion(true) ? 0 : -1;
    }

    public static void main(String[] args) throws Exception {
        if (args.length != 2) {
            System.err.println("Usage: WordCount <in> <frontiers> <out>");
            System.exit(2);
        }
        if (run(args[0], args[1]) < 0) {
            System.exit(1);
        }
    }
}
```

6: Spark Programming

PageRank:

```
import sys
from pyspark import SparkConf, SparkContext
import time

if __name__ == '__main__':
    conf = SparkConf()
    sc = SparkContext(conf=conf)
    sc.setLogLevel("ERROR") # Added to avoid Warnings cluttering the terminal
    file_path = sys.argv[1]
    lines = sc.textFile(file_path)

    # Parameters
    damping = 0.8
    num_iterations = 50
    top_i_nodes = 5
    dynamic = sys.argv[2]

    first = time.time()

    # Parse the lines into (source, destination) pairs and remove duplicates
    edges = lines.map(lambda line: tuple(map(int, line.split()))).distinct()

    # We estimate the amount of nodes
    if(dynamic):
        # This method takes aprox. 0.5s longer but is dynamic
        nodes = edges.flatMap(lambda edge: edge).distinct()
        n = nodes.count()
    else:
        # This method is faster but not dynamic
        n = 100 if "small" in file_path else 1000 if "full" in file_path else None

    # Create an adjacency list as (node, [neighbors])
    adj_list = edges.groupByKey().mapValues(list).cache()

    # Initialize each node's PageRank value
    page_ranks = adj_list.mapValues(lambda _: 1.0 / n)

    for i in range(num_iterations):
        # Broadcast the adjacency list for efficient access
        adjacency_broadcast = sc.broadcast(adj_list.collectAsMap())

        # Compute contributions for each node's neighbors
        contributions = page_ranks.flatMap(lambda node_rank: [
            (neighbor, node_rank[1] /
            len(adjacency_broadcast.value.get(node_rank[0], [])))
            for neighbor in adjacency_broadcast.value.get(node_rank[0], [])
        ])

        # Aggregate contributions and calculate new PageRank values
        page_ranks = contributions.reduceByKey(lambda a, b: a + b).mapValues(
            lambda rank: (1 - damping) / n + damping * rank
        )

    # Get the top 5 nodes with the highest PageRank scores
    highest = page_ranks.takeOrdered(top_i_nodes, key=lambda x: -x[1])

    # Print the top 5 nodes
    print("Top 5 nodes with highest PageRank scores:")
    for node, score in highest:
        print(f"Node {node}: {score}")

    last = time.time()
    print("Total program time: %.2f seconds" % (last - first))
```

```
sc.stop()
```

WordCount:

```
import re
import sys
from pyspark import SparkConf, SparkContext
import time

if __name__ == '__main__':
    conf = SparkConf()
    sc = SparkContext(conf=conf)
    sc.setLogLevel("ERROR") # Added to avoid Warnings cluttering the terminal
    lines = sc.textFile(sys.argv[1])
    top_i_words = 10

    first = time.time()

    # We split the lines into words
    words = lines.flatMap(lambda line: re.split(r'[^\w]+', line))

    # We now count every word
    word_counts = words.countByValue()
    top_words = sorted(word_counts.items(), key=lambda x: -x[1])[:top_i_words]

    # Print the top i words
    print(f"Top {top_i_words} words:")
    for word, count in top_words:
        print(f"({repr(word)}, {count})")

    last = time.time()

    print("Total program time: %.2f seconds" % (last - first))
    sc.stop()
```

7: Worse is better

What I learned

I find the article “Worse is Better” by Richard P. Gabriel interesting as, while somewhat dated, it gives an interesting view into the compromises and decisions made in the early days of programming and computer science.

The article explores two contrasting design philosophies in software development: the MIT approach, which focuses on achieving a complete and ideal “right” design, and the “worse-is-better” philosophy, also called the New Jersey approach, which values simplicity and practicality. Gabriel argues that simpler, less complete designs (termed “worse”) are often more successful in real-world applications, as they foster adaptability and ease of use. This approach prioritizes implementation simplicity over full correctness or completeness, leading to systems that are easier to port and modify. He illustrates this with Unix and C, where simplicity and minimal resource requirements enable these systems to thrive across various environments. The title, **The Rise of Worse is Better**, reflects how “worse” design choices—those sacrificing some ideal goals—often yield “better” real-world success through increased portability, adaptability, and longevity. In this context, “worse” signifies compromises in design ideals, while “better” represents the widespread adoption and durability achieved by these practical systems.

8: Graph Partitioning

GraphPartitioner.py:

```
import sys
import struct
import random
from collections import defaultdict
import matplotlib.pyplot as plt
import networkx as nx

class GraphPartitioner:
    def __init__(self, file_path, num_partitions, show_details):
        self.file_path = file_path
        self.num_partitions = num_partitions
        self.show_details = show_details
        self.edges = self.read_graph()
        self.directed = self.isDirected()

    def isDirected(self):
        # Some of the graphs are undirected
        return not any(file_name in file_path for file_name in ["small-5",
"synthesized-1b"])

    # A (Slightly modified) replication of the C code already existing
    # Reads the data and saves all edges in the graph
    def read_graph(self):
        edges = []
        with open(self.file_path, "rb") as file:
            data = file.read(8)      # Read first 8 bytes
            while data:
                src, dst = struct.unpack("ii", data) # (4 byte src, 4 byte dst)
                edges.append((src, dst))
                data = file.read(8) # Read the next 8 bytes
        print("Done reading edges\n")
        return edges

    # Show graphs in subplots for each partition
    def show_graph(self, partitions, title="Graph Partitions"):
        height = 4
        width = len(partitions) * height
        fig, axes = plt.subplots(1, len(partitions), figsize=(width, height))

        if len(partitions) == 1:
            axes = [axes] # Ensure axes is always iterable for a single partition

        for i, (partition_id, partition) in enumerate(sorted(partitions.items())):
            ax = axes[i]
            G = nx.DiGraph() if self.directed else nx.Graph()
            G.add_edges_from(partition['edges'])

            # Differentiate master vertices by color, or color all nodes lightgreen
            if none exist
                master_vertices = partition.get('master_vertices', set())
                # If there is no master_vertices value, then we are on the Full graph
                and all are master vertices
                if not master_vertices:
                    node_colors = ["lightgreen" for _ in G.nodes()]
                    # If only some are master vertices, color them lightgreen and the rest
                    skyblue
                else:
                    node_colors = ["lightgreen" if node in master_vertices else
"skyblue" for node in G.nodes()]

            pos = nx.spring_layout(G) # Use spring layout for better visualization
            nx.draw(G, pos, with_labels=True, node_color=node_colors,
font_weight="bold",
```

```

        node_size=500, edge_color="gray", ax=ax, arrows=self.directed)
ax.set_title(f"Partition {partition_id}")

fig.suptitle(title, fontsize=16)
plt.tight_layout(rect=[0, 0, 1, 0.95])
plt.show()

# Balanced p-way edge-cut partitioning
def edge_cut_partition(self):
    partitions = defaultdict(
        lambda: {
            'master_vertices': set(),
            'total_vertices': set(),
            'replicated_edges': 0,
            'edges': []
        }
    )
    vertex_to_partition = {}

    for src, dst in self.edges:
        # Hash vertices to assign them to partitions
        src_partition = vertex_to_partition.get(src, hash(src) %
self.num_partitions)
        dst_partition = vertex_to_partition.get(dst, hash(dst) %
self.num_partitions)

        # Assign the vertex to the partition
        vertex_to_partition[src] = src_partition
        vertex_to_partition[dst] = dst_partition

        # Assign edges and count replicated edges
        if src_partition != dst_partition:
            partitions[src_partition]['replicated_edges'] += 1
            partitions[dst_partition]['replicated_edges'] += 1

        partitions[src_partition]['edges'].append((src, dst))
        partitions[dst_partition]['edges'].append((src, dst))

        # Update master and total vertices
        partitions[src_partition]['master_vertices'].add(src)
        partitions[dst_partition]['master_vertices'].add(dst)
        partitions[src_partition]['total_vertices'].update([src, dst])
        partitions[dst_partition]['total_vertices'].update([src, dst])

    # Print partition statistics
    for i, partition in sorted(partitions.items()):
        print(f"Partition {i}")
        print(len(partition['master_vertices']))
        print(len(partition['total_vertices']))
        print(partition['replicated_edges'])
        print(len(partition['edges']))

    # Show the partitions side by side
    if self.show_details:
        self.show_graph(partitions, title="Edge-Cut Partitioning")

# Balanced p-way vertex-cut partitioning
def vertex_cut_partition(self):
    partitions = defaultdict(
        lambda: {
            'master_vertices': set(),
            'total_vertices': set(),
            'edges': []
        }
    )
    vertex_partitions = defaultdict(set) # Store all partitions a vertex is
assigned to

```

```

for src, dst in self.edges:
    # Hash the edge to assign it to a partition
    edge_partition = hash((src, dst)) % self.num_partitions
    partitions[edge_partition]['edges'].append((src, dst))

    # Update vertex partitions
    vertex_partitions[src].add(edge_partition)
    vertex_partitions[dst].add(edge_partition)

# Calculate master and total vertices for each partition
for vertex, assigned_partitions in vertex_partitions.items():
    # Choose one partition as master
    master_partition = random.choice(list(assigned_partitions))
    for partition_id in assigned_partitions:
        partitions[partition_id]['total_vertices'].add(vertex)
        if partition_id == master_partition:
            partitions[partition_id]['master_vertices'].add(vertex)

# Print partition statistics
for i, partition in sorted(partitions.items()):
    print(f"Partition {i}")
    print(f"{len(partition['master_vertices'])}")
    print(f"{len(partition['total_vertices'])}")
    print(f"{len(partition['edges'])}")

# Show the partitions side by side
if self.show_details:
    self.show_graph(partitions, title="Vertex-Cut Partitioning")

# Greedy heuristic vertex-cut partitioning
def greedy_heuristic_partition(self):
    partitions = defaultdict(
        lambda: {
            'master_vertices': set(),
            'total_vertices': set(),
            'edges': []
        }
    )
    vertex_degrees = defaultdict(int)

    # Calculate degrees for all vertices
    for src, dst in self.edges:
        vertex_degrees[src] += 1
        vertex_degrees[dst] += 1

    # Assign vertices to partitions using greedy heuristic
    for src, dst in self.edges:
        src_partition = vertex_degrees[src] % self.num_partitions
        dst_partition = vertex_degrees[dst] % self.num_partitions

        partitions[src_partition]['edges'].append((src, dst))
        partitions[dst_partition]['edges'].append((src, dst))

        partitions[src_partition]['total_vertices'].add(src)
        partitions[dst_partition]['total_vertices'].add(dst)

        master_partition = random.choice([src_partition, dst_partition])
        if master_partition == src_partition:
            partitions[src_partition]['master_vertices'].add(src)
        else:
            partitions[dst_partition]['master_vertices'].add(dst)

    # Ensure every partition gets at least one vertex
    for i in range(self.num_partitions):
        if not partitions[i]['total_vertices']:
            random_vertex = random.choice(list(vertex_degrees.keys()))
            partitions[i]['total_vertices'].add(random_vertex)
            partitions[i]['master_vertices'].add(random_vertex)

```

```

# Print partition statistics
for i, partition in sorted(partitions.items()):
    print(f"Partition {i}")
    print(len(partition['master_vertices']))
    print(len(partition['total_vertices']))
    print(len(partition['edges']))

# Show the partitions side by side
if self.show_details:
    self.show_graph(partitions, title="Greedy Vertex-Cut Partitioning")

# Helper function to get the neighbors of a vertex
def get_neighbors(self, vertex):
    neighbors = set()
    for src, dst in self.edges:
        if src == vertex:
            neighbors.add(dst)
        elif dst == vertex:
            neighbors.add(src)
    return neighbors

# Helper function to get partition for a vertex
def get_partition_for_vertex(self, vertex, partitions):
    for partition_id, partition in partitions.items():
        if vertex in partition['total_vertices']:
            return partition_id
    return None

# Balanced p-way Hybrid-Cut based on PowerLyra
def hybrid_cut_partition(self, theta=10):
    partitions = defaultdict(
        lambda: {
            'master_vertices': set(),
            'total_vertices': set(),
            'edges': []
        }
    )
    vertex_degree = defaultdict(int) # Store vertex degrees
    vertex_partitions = defaultdict(set) # Tracks the partitions each vertex
is assigned to

    # Calculate degrees for all vertices
    for src, dst in self.edges:
        vertex_degree[src] += 1
        vertex_degree[dst] += 1

    for src, dst in self.edges:
        if vertex_degree[src] > theta or vertex_degree[dst] > theta:
            # High-degree vertices: use edge-cut strategy
            src_partition = hash(src) % self.num_partitions
            dst_partition = hash(dst) % self.num_partitions
            partitions[src_partition]['edges'].append((src, dst))
            partitions[dst_partition]['edges'].append((src, dst))
            partitions[src_partition]['master_vertices'].add(src)
            partitions[dst_partition]['master_vertices'].add(dst)
            partitions[src_partition]['total_vertices'].update([src, dst])
            partitions[dst_partition]['total_vertices'].update([src, dst])
        else:
            # Low-degree vertices: use vertex-cut (greedy heuristic)
            min_partition = None
            min_replication = float('inf')

            for partition_id in range(self.num_partitions):
                replication_cost = (
                    int(src not in partitions[partition_id]['master_vertices'])
+
                    int(dst not in partitions[partition_id]['master_vertices']))

```

```

        )
    if replication_cost < min_replication:
        min_replication = replication_cost
        min_partition = partition_id

    partitions[min_partition]['edges'].append((src, dst))
    partitions[min_partition]['master_vertices'].update([src, dst])
    partitions[min_partition]['total_vertices'].update([src, dst])
    vertex_partitions[src].add(min_partition)
    vertex_partitions[dst].add(min_partition)

# Print partition statistics
for i, partition in sorted(partitions.items()):
    print(f"Partition {i}")
    print(f"{len(partition['master_vertices'])}"))
    print(f"{len(partition['total_vertices'])}"))
    print(f"{len(partition['edges'])}"))

if self.show_details:
    self.show_graph(partitions, title=f"Hybrid-Cut Partitioning
(Theta={theta})")

if __name__ == "__main__":
    file_path = sys.argv[1] if len(sys.argv) > 1 else "small-5.graph"
    num_partitions = int(sys.argv[2]) if len(sys.argv) > 2 else 3
    show_details = (sys.argv[3]!="False") if len(sys.argv) > 3 else False
    partitioner = GraphPartitioner(file_path, num_partitions, show_details)

    # Display the graph if applicable
    if show_details:
        partitioner.show_graph({0: {'edges': partitioner.edges}}, title="Full
Graph")

    # Partition the graph
    print("\nEdge-Cut Partitioning:")
    partitioner.edge_cut_partition()
    print("\nVertex-Cut Partitioning:")
    partitioner.vertex_cut_partition()
    print("\nHybrid-Cut Partitioning:")
    partitioner.hybrid_cut_partition(theta=2)
    print("\nGreedy-Cut Partitioning:")
    partitioner.greedy_heuristic_partition()

```

9: Graph algorithms with GridGraph

Kcores.cpp:

```
#include "core/graph.hpp"

int main(int argc, char **argv) {
    if (argc < 3) {
        fprintf(stderr, "usage: kcores [path] [k] [memory budget in GB]\n");
        exit(-1);
    }

    std::string path = argv[1];
    int k = atoi(argv[2]);
    long memory_bytes = (argc >= 4) ? atol(argv[3]) * 10241 * 10241 * 10241 : 81 * 10241 * 10241 * 10241;

    Graph graph(path);
    graph.set_memory_bytes(memory_bytes);

    Bitmap *active_in = graph.alloc_bitmap();
    Bitmap *active_out = graph.alloc_bitmap();
    BigVector<int> degree(graph.path + "/degree", graph.vertices);
    BigVector<int> core(graph.path + "/core", graph.vertices);

    long vertex_data_bytes = (long)graph.vertices * (sizeof(int) + sizeof(int));
    graph.set_vertex_data_bytes(vertex_data_bytes);

    // Initialize degree and active vertices
    active_out->fill();
    degree.fill(0);
    graph.stream_edges<VertexId>(
        [&](Edge &e) {
            write_add(&degree[e.source], 1);
            return 0;
        },
        nullptr, 0, 0);

    // Initialize core and set initial active vertices
    int active_vertices = graph.stream_vertices<VertexId>(
        [&](VertexId i) {
            core[i] = (degree[i] >= k) ? 1 : 0;
            return core[i];
        });
}

printf("Initialization complete: %d active vertices\n", active_vertices);

// K-core decomposition iterations
int iteration = 0;
while (active_vertices > 0) {
    iteration++;
    printf("Iteration %d: %d active vertices\n", iteration, active_vertices);

    std::swap(active_in, active_out);
    active_out->clear();

    graph.hint(degree, core);
    active_vertices = graph.stream_edges<VertexId>(
        [&](Edge &e) {
            if (core[e.source] == 1 && core[e.target] == 0) {
                write_add(&degree[e.target], -1);
                if (degree[e.target] < k) {
                    core[e.target] = 0;
                    active_out->set_bit(e.target);
                    return 1;
                }
            }
        })
}
```

```

        return 0;
    },
    active_in);
}

// Count k-core vertices
int kcore_vertices = graph.stream_vertices<VertexId>(
    [&](VertexId i) {
        return core[i] == 1;
});

printf("K-core (%d-core) decomposition complete: %d vertices remain\n", k,
kcore_vertices);

return 0;
}

```

PageRankDelta.cpp:

```

#include "core/graph.hpp"

int main(int argc, char ** argv) {
    if (argc < 3) {
        fprintf(stderr, "usage: pagerank_delta [path] [iterations] [memory
budget in GB]\n");
        exit(-1);
    }
    std::string path = argv[1];
    int iterations = atoi(argv[2]);
    long memory_bytes = (argc >= 4) ? atol(argv[3]) * 1024l * 1024l * 1024l : 81
* 1024l * 1024l * 1024l;

    Graph graph(path);
    graph.set_memory_bytes(memory_bytes);

    BigVector<VertexId> degree(graph.path + "/degree", graph.vertices);
    BigVector<float> pagerank(graph.path + "/pagerank", graph.vertices);
    BigVector<float> delta(graph.path + "/delta", graph.vertices);
    BigVector<float> new_delta(graph.path + "/new_delta", graph.vertices);

    long vertex_data_bytes = (long)graph.vertices * (sizeof(VertexId) +
sizeof(float) * 3);
    graph.set_vertex_data_bytes(vertex_data_bytes);

    double begin_time = get_time();

    // Initialize degrees
    degree.fill(0);
    graph.stream_edges<VertexId>(
        [&](Edge & e) {
            write_add(&degree[e.source], 1);
            return 0;
        }, nullptr, 0, 0
    );

    // Initialize Pagerank and Delta
    graph.hint(pagerank, delta, new_delta);
    graph.stream_vertices<VertexId>(
        [&](VertexId i) {
            pagerank[i] = 0.15f; // Initial PageRank value
            delta[i] = 1.0f / degree[i]; // Initial delta
            new_delta[i] = 0;
            return 0;
        }, nullptr, 0,
        [&](std::pair<VertexId, VertexId> vid_range) {
            pagerank.load(vid_range.first, vid_range.second);
            delta.load(vid_range.first, vid_range.second);
            new_delta.load(vid_range.first, vid_range.second);
        }
    );
}

```

```

        new_delta.load(vid_range.first, vid_range.second);
    },
    [&] (std::pair<VertexId, VertexId> vid_range) {
        pagerank.save();
        delta.save();
        new_delta.save();
    }
);

// PageRank Delta Iterations
for (int iter = 0; iter < iterations; iter++) {
    graph.hint(delta, new_delta);
    graph.stream_edges<VertexId>(
        [&] (Edge & e) {
            write_add(&new_delta[e.target], 0.85f * delta[e.source]);
            return 0;
        }, nullptr, 0, 1,
        [&] (std::pair<VertexId, VertexId> source_vid_range) {
            delta.lock(source_vid_range.first,
source_vid_range.second);
        },
        [&] (std::pair<VertexId, VertexId> source_vid_range) {
            delta.unlock(source_vid_range.first,
source_vid_range.second);
        }
    );
    graph.hint(pagerank, delta, new_delta);
    graph.stream_vertices<float>(
        [&] (VertexId i) {
            pagerank[i] += new_delta[i];
            delta[i] = new_delta[i] / degree[i];
            new_delta[i] = 0; // Reset for next iteration
            return 0;
        }, nullptr, 0,
        [&] (std::pair<VertexId, VertexId> vid_range) {
            pagerank.load(vid_range.first, vid_range.second);
            delta.load(vid_range.first, vid_range.second);
            new_delta.load(vid_range.first, vid_range.second);
        },
        [&] (std::pair<VertexId, VertexId> vid_range) {
            pagerank.save();
            delta.save();
            new_delta.save();
        }
    );
}

double end_time = get_time();
printf("%d iterations of pagerank delta took %.2f seconds\n", iterations,
end_time - begin_time);

return 0;
}

```

10: Graph Mining

FSM.py

```
import pandas as pd
import numpy as np
import json
from collections import defaultdict
from itertools import combinations
from networkx.algorithms import isomorphism
from math import comb
from multiprocessing import Pool

def build_graph(edges, vertices):
    print("Building graph...")
    graph = defaultdict(list)

    # Ensure unique indices in the vertices DataFrame
    vertices = vertices.drop_duplicates(subset='id').set_index('id')

    # Add all edges to the graph
    for _, row in edges.iterrows():
        source, target = row['source_id'], row['target_id']
        graph[source].append((target, row['amt'], row['strategy_name'],
        row['buscode']))

    print(f"Total nodes: {len(graph)}, Total edges: {sum(len(v) for v in
graph.values())}\n")
    return graph

def hash_edge(source, target, amt, strategy_name, buscode):
    # I think we can go without some of these values, but thought they might be
    nice to have anyways :-)
    return f"{min(source, target)}-{max(source, target)}-{amt}-{strategy_name}-
{buscode}"

def mine_frequent_subgraphs(graph, pattern_size, support_threshold, output_file):
    print("Mining frequent subgraphs...")
    subgraph_counts = defaultdict(int)
    subgraph_patterns = defaultdict(list)

    # Iterate over node combinations
    for nodes in combinations(graph.keys(), pattern_size):
        edges = []
        for u, v in combinations(nodes, 2):
            if v in [neighbor[0] for neighbor in graph[u]]:
                edge_details = next(
                    (neighbor for neighbor in graph[u] if neighbor[0] == v), None
                )
                if edge_details:
                    edges.append(
                        hash_edge(u, v, *edge_details[1:])
                    )

        if len(edges) == pattern_size:
            subgraph_key = "_".join(sorted(edges))
            subgraph_counts[subgraph_key] += 1
            subgraph_patterns[subgraph_key].append(edges)

    # Filter frequent subgraphs
    frequent_subgraphs = {
        k: v for k, v in subgraph_counts.items() if v >= support_threshold
    }

    save_results(frequent_subgraphs, subgraph_patterns, output_file)
    print("Frequent subgraph mining completed.")
```

```

def save_results(frequent_subgraphs, subgraph_patterns, output_file):
    result = []
    for subgraph, frequency in frequent_subgraphs.items():
        edges = subgraph_patterns[subgraph]
        result.append({
            "frequency": frequency,
            "edges": [{"source": edge.split("-") [0], "target": edge.split("-") [1],
            "details": edge} for edge in edges]
        })

    with open(output_file, 'w') as f:
        json.dump(result, f, indent=4)
    print(f"Results saved to {output_file}")

if __name__ == "__main__":
    print("Reading data...")
    header1 = ['id', 'name', 'timestamp', 'black']
    header2 = ['source_id', 'target_id', 'timestamp', 'amt', 'strategy_name',
    'trade_no', 'buscode', 'other']

    account = pd.read_csv('data/account', names=header1, sep=',')
    card = pd.read_csv('data/card', names=header1, sep=',')
    account_to_account = pd.read_csv('data/account_to_account', names=header2,
    sep=',', usecols=range(len(header2)))
    account_to_card = pd.read_csv('data/account_to_card', names=header2, sep=',',
    usecols=range(len(header2)))

    vertices = pd.concat([account, card])
    edges = pd.concat([account_to_account, account_to_card])
    edges['amt'] = edges['amt'].round()

    graph = build_graph(edges, vertices)

    # Start mining frequent subgraphs
    mine_frequent_subgraphs(
        graph,
        pattern_size=3,
        support_threshold=10000,
        output_file=f"results/bdci_data.json"
    )

```

11: Spark Streaming Top-K

Top-k.py:

```
from __future__ import print_function
import sys
from pyspark import SparkContext
from pyspark.streaming import StreamingContext
from pyspark.streaming.dstream import DStream

# --- Code Setup ---

# Initialize SparkContext and StreamingContext
sc = SparkContext(appName="Py_HDFSWordCount")
ssc = StreamingContext(sc, 60)

# Create a DStream that listens to the HDFS directory
hdfs_directory = "hdfs://intro00:8020/user/2024403421/stream"
lines = ssc.textFileStream(hdfs_directory)

# We use these global variables to keep track of our top-k algorithm
word_counts = {}
file_no = 1
last_file = 5
k = 100

# --- My two functions to perform the Top-K ---

# Function to update the word counts with each new RDD
def update_count(new_counts, last_counts):
    # On first batch we initialize an empty dictionary
    if last_counts is None:
        last_counts = {}

    # On all continuing files, update the word counts
    for word, count in new_counts:
        if word in last_counts:
            last_counts[word] += count
        else:
            last_counts[word] = count

    return last_counts

# Function to process each RDD and compute the top-k frequent words
def process_rdd(time, rdd):
    global word_counts, file_no, last_file, k
    if rdd.isEmpty():
        return

    # Compute word counts for the current RDD
    counts = rdd.flatMap(lambda line: line.split(" ")) \
        .map(lambda word: (word, 1)) \
        .reduceByKey(lambda a, b: a + b)

    # Update the global word counts using the updateCounts function
    updated_counts = counts.collect() # Collect current counts to update
    word_counts = update_count(updated_counts, word_counts)

    # Sort by count and take top k
    top_k = sorted(word_counts.items(), key=lambda x: x[1], reverse=True)[:k]

    # We structure the output
    output = f"---- File Number {file_no} ----\n"
    output += f" --- Top-{k} words so far ---\n"

    for word, count in top_k:
        output += f"{word}: {count}\n"

    file_no += 1
    last_file = file_no
```

```
# We write the output to a file
output_filename = f"output_file_{file_no}.txt"
with open(output_filename, "w") as f:
    f.write(output)

print(output)

file_no += 1

# Stop the program after processing the last file
if file_no > last_file:
    print("\nMaximum number of files processed. Stopping the streaming.")
    ssc.stop(stopSparkContext=True, stopGraceFully=True)

# --- Use My functions ---

# Process each RDD in the DStream and compute top-k
lines.foreachRDD(process_rdd)

# Start streaming and wait for termination
ssc.start()
ssc.awaitTermination()
```

12: Halide Scheduling

Dilated-conv.cpp:

```
#include "Halide.h"
#include "common.h"

#include <stdio.h>

using namespace Halide;
using namespace Halide::Tools;

int main(int argc, char **argv) {
    const int N = 5, CI = 128, CO = 128, W = 100, H = 80, KW = 3, KH = 3;
    const int dilation = 15;

    ImageParam input(type_of<float>(), 4);
    ImageParam filter(type_of<float>(), 4);

    // Define variables and reduction domain
    Var x("x"), y("y"), c("c"), n("n");
    Var xo("xo"), yo("yo"), xi("xi"), yi("yi");
    Var co("co"), ci("ci");

    Func dilated_conv("dilated_conv");
    RDom r(0, CI, 0, KW, 0, KH);

    // Algorithm definition
    dilated_conv(c, x, y, n) = 0.0f;
    dilated_conv(c, x, y, n) += filter(c, r.y, r.z, r.x) *
        input(r.x, x + r.y * (dilation + 1), y + r.z * (dilation + 1), n);

    // **Scheduling**
    // 1. Split the x and y dimensions for tiling
    dilated_conv.compute_root()
        .tile(x, y, xo, yo, xi, yi, 8, 8) // 8x8 tile size (tunable)
        .fuse(xo, yo, co) // Parallelize outer loop over tiles
        .parallel(co) // Vectorize inner x dimension

    // 2. Optimize the reduction
    dilated_conv.update()
        .reorder(r.x, r.y, r.z, c, xi, yi, n)
        .unroll(r.z) // Unroll kernel height loop
        .unroll(r.y) // Unroll kernel width loop
        .vectorize(xi, 8) // Vectorize inner computation
        .parallel(n); // Parallelize across batch dimension

    // Buffer initialization
    Buffer<float, 4> in(CI, W + (KW - 1) * (dilation + 1), H + (KH - 1) * (dilation + 1), N);
    Buffer<float, 4> fil(CO, KW, KH, CI);
    Buffer<float, 4> output_halide(CO, W, H, N);

    // Initialize input and filter with random data
    random_data<float, 4>(in);
    random_data<float, 4>(fil);
    input.set(in);
    filter.set(fil);

    // JIT compile and run
    dilated_conv.realize(output_halide);
    double t_halide = benchmark(10, 10, [&]()
    { dilated_conv.realize(output_halide); });

    Buffer<float, 4> output_ref(CO, W, H, N);
```

```

    double t_onednn = dnnl_dilated_conv_wrapper(in.data(), fil.data(),
output_ref.data(),
                                         {N, CI, CO, W, H, KW, KH, dilation,
dilation});

    // Check correctness
    if (check_equal<float, 4>(output_ref, output_halide)) {
        printf("Halide results - OK\n");
    } else {
        printf("Halide results - FAIL\n");
        return 1;
    }

    float gflops = 2.0f * (N * CO * H * W) * (CI * KH * KW) / 1e9f;
    printf("Halide: %fms, %f GFLOP/s\n", t_halide * 1e3, (gflops / t_halide));
    printf("oneDNN: %fms, %f GFLOP/s\n\n", t_onednn * 1e3, (gflops / t_onednn));

    printf("Success!\n");
    return 0;
}

```

| | |
|---|-----------|
| BDS Exam Preparations..... | 2 |
| Parallel Programming..... | 2 |
| Amdahl's Law..... | 2 |
| Gustafson's Law..... | 2 |
| Realistic Data Pipeline..... | 5 |
| Estimating Power Efficiency of Dual Core CPU..... | 5 |
| OpenMP..... | 5 |
| Data Sharing:..... | 6 |
| MPI Programs..... | 7 |
| Parallel I/O Systems..... | 9 |
| Fault Tolerance..... | 9 |
| Disk Level..... | 9 |
| RAID..... | 9 |
| Node Level..... | 10 |
| Google File System (GFS)..... | 11 |
| MapReduce..... | 12 |
| In Memory Computing..... | 14 |
| RDD (Resilient Distributed Datasets)..... | 15 |
| Spark Programming..... | 15 |
| Graph Analysis..... | 17 |
| Asynchronous processing..... | 17 |
| Graph partitioning on power law graphs and GAS model..... | 17 |
| Out of core graph processing..... | 18 |
| GRAPHCHI – Parallel Sliding Window..... | 19 |
| XSTREAM – Edge Centric..... | 19 |
| Streaming Partitions..... | 19 |
| GRIDGRAPH – 2D Partition..... | 20 |
| Top Down vs. Bottom Up..... | 20 |
| GEMINI – Distributed push/pull..... | 21 |
| Shentu - Designed for Extreme Scale Graph..... | 21 |
| Graph Mining..... | 22 |
| Stream Processing..... | 22 |
| Actor Model..... | 22 |
| Stream Processing – STORM..... | 23 |
| Core Concepts and Architecture..... | 23 |
| SparkStream..... | 24 |
| Comparison of different Streaming Models..... | 24 |
| Machine Learning..... | 24 |

| | |
|---|----|
| Halide..... | 25 |
| Graph Optimization..... | 25 |
| Distributed Training..... | 25 |
| Memory Consumption..... | 26 |
| AdaPipe..... | 26 |
| Blockchain..... | 27 |
| Conflux..... | 28 |
| Computer Systems..... | 29 |
| Why is it difficult to build a system?..... | 29 |
| Common Problems of Systems..... | 29 |
| Coping with Complexity..... | 29 |
| Prep. Exam..... | 30 |

BDS Exam Preparations

Challenges of Big Data Processing:

- Volume → More memory and I/O bandwidth
- Velocity → Fast Processing
- Variety → More computing paradigms

Parallel Programming

Challenges:

- Parallelism identification
- Load balancing
- Proper synchronization

Amdahl's Law

If P is the proportion of a program that can be made parallel, the limit of speedup of whole program is bounded at **Speed-up = 1/(1-P)**

Gustafson's Law

Amdahl's law assumes fixed problem size

If problem size can increase as the compute power increases, the speedup that can be achieved is much larger.

$$S = a(n) + p(1-a(n))$$

n = problem size

p = number of processes

a(n) = serial portion at problem size n

So $S \rightarrow p$ when $a(n) = 0$

⇒ Larger machines are still useful

⇒ Many small chips can deliver good performance when the problem size is big

Problem: Some applications don't have bigger data set

Real world is mix of Amdahl's law and Gustafson's law

Posix Threads (Pthreads) Interface

-

Pthreads:

Standard interface for ~60 functions that

manipulate threads from C programs

–Creating and reaping threads

- `pthread_create()`

- `pthread_join()`

–Determining your thread ID

- `pthread_self()`

–Terminating threads

- `pthread_cancel()`

- `pthread_exit()`

• `exit()` [terminates all threads] ,`RET` [terminates current thread]

–Synchronizing access to shared variables

- `pthread_mutex_init`

- `pthread_mutex_[un]lock`

- `pthread_cond_init`

- `pthread_cond_[timed]wait`

Semaphores:

Non-negative global integer synchronization variables

Manipulated by P and V operations:

- $P(s)$: [`while (s == 0) wait(); s--;`]

- Dutch for "Proberen" (test)

- $V(s)$: [`s++;`]

- Dutch for "Verhogen" (increment)

OS kernel guarantees that operations between brackets [] are executed indivisibly

- Only one P or V operation at a time can modify s.
- When while loop in P terminates, only that P can decrement s

```
void *sum_thread(void *vargp)
{
    long sum;
    int myid = (long) vargp;
    int start = myid * sample_num / 4;
    int end = (myid+1) * sample_num / 4 - 1;

    sum = 0;
    for (int i=start; i<=end; i++)
        sum += x[i];

    sem_wait(&num_mutex);
    global_sum += sum;
    sem_post(&num_mutex);

    return NULL;
}
```

Issues with Pthreads:

Complex for data processing

- Pthread is not always portable(e.g. Windows)
- Code to calculate addresses and loop bounds
- May not be easy to specify scheduling schemes

```
void *sum_thread(void *);

volatile long global_sum = 0; /* global */

main()
{
    ...
    for (int i = 0; i < N; i++)
        pthread_create(&tid[i], NULL, sum_thread, (void*)i);

    for (int i = 0; i < N; i++)
        pthread_join(tid[i], NULL);
}
```

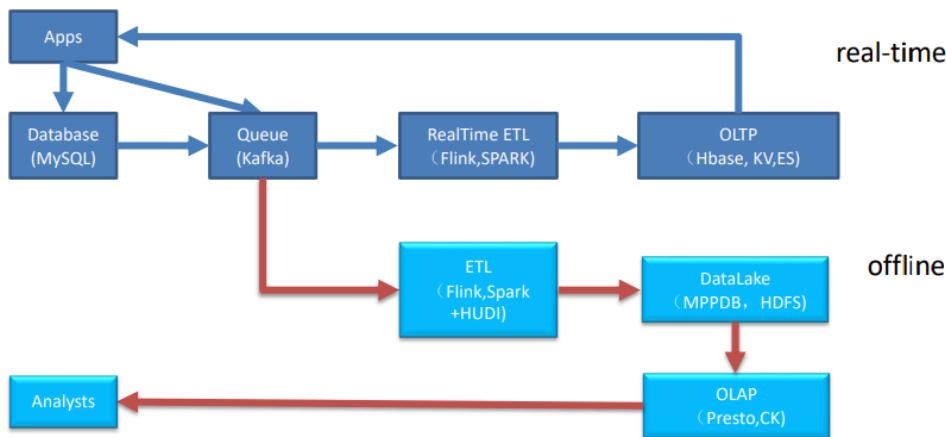
```
void *sum_thread(void *);

volatile long global_sum = 0; /* global */

main()
{
    ...
    for (int i = 0; i < N; i++)
        pthread_create(&tid[i], NULL, sum_thread, (void*)i);

    for (int i = 0; i < N; i++)
        pthread_join(tid[i], NULL);
}
```

Realistic Data Pipeline



Estimating Power Efficiency of Dual Core CPU

Single Core: 1 solution/s

$$P = cf_1^2$$

$$P/2 = cf_2^2$$

$$F_2 = f_1/(2)^{1/2}$$

$$\text{Dual Core Throughput} = 2*f_2 = 2^{1/2} * f_1 = 2^{1/2} \text{ solutions/s}$$

⇒ With perfect parallelism

OpenMP

“Standard” API for defining multi-threaded sharedmemory programs

– Allow a programmer to **separate a program into serial regions and parallel regions**, rather than T concurrently-executing threads.

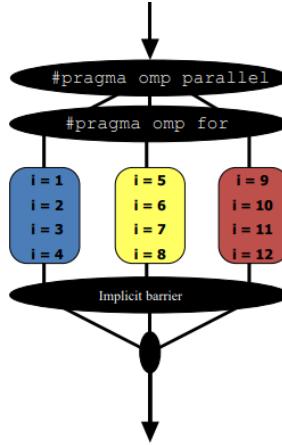
– **Hide stack management**

– Provide synchronization constructs

⇒ **ONLY SINGLE MACHINES**

- Threads are assigned an independent set of iterations
- Threads must wait at the end of work-sharing construct

```
#pragma omp parallel for
for(i = 1, i < 13, i++)
    c[i] = a[i] + b[i]
```



Data Sharing:

- Parallel programs often employ two types of data
 - Shared data, visible to all threads, similarly named
 - Private data, visible to a single thread (often stack-allocated)
- pthreads:
 - Global-scoped variables are shared
 - Stack-allocated variables are private
- OpenMP:
 - **shared** variables are shared
 - **private** variables are private
 - Rules for default

```
int bigdata[1024];

void* foo(void* bar) {
    int tid;

    #pragma omp parallel \
        shared ( bigdata ) \
        private ( tid )
    {
        /* Calc. here */
    }
}

int i = 0;
int n = 10;
int a = 7;

#pragma omp parallel for
for (i = 0; i < n; i++)
{
    int b = a + i;
    ...
}
```

n and a are shared. i and b are private.

Firstprivate

```
#include <omp.h>
#include <stdio.h>

int main()
{
    int sum = 0;

    #pragma omp parallel for firstprivate(sum)
    for (int i = 1; i <= 10; i++)
    {
        sum += i;
        printf("In iteration %d, sum is %d\n", i, sum);
    }

    printf("After the loop, sum is %d\n", sum);
}
```

- Initialize private variables from the master threads

lastprivate

```
#include <omp.h>
#include <stdio.h>

int main()
{
    int sum = 0;

    #pragma omp parallel for firstprivate(sum) lastprivate(sum)
    for (int i = 1; i <= 10; i++)
    {
        sum += i;
        printf("In iteration %d, sum is %d\n", i, sum);
    }

    printf("After the loop, sum is %d\n", sum);
}
```

- Pass the value of a private from the **last iteration** to global variable

OpenMP Critical Construct

- `#pragma omp critical [(lock_name)]`
- Defines a critical region on a structured block

```
Threads wait their turn -at
a time, only one calls
consum() thereby
protecting RES from race
conditions
```

```
float RES;
#pragma omp parallel
{ float B;
#pragma omp for private(B)
for(int i=0; i<niters; i++){
    B = big_job(i);
#pragma omp critical
    {
        consum (B, RES);
    }
}
```

Sum of an Array

```
#include <omp.h>
#include <stdio.h>

int main()
{
    int sum = 0;

#pragma omp parallel for reduction(+:sum)
for (int i = 1; i <= 10; i++)
{
    sum += i;
    printf("In iteration %d, sum is %d\n", i, sum);
}

printf("After the loop, sum is %d\n", sum);
```

- The `schedule clause` affects how loop iterations are mapped onto threads

```
schedule (static [,chunk])
    • Blocks of iterations of size "chunk" to threads
    • Round robin distribution
    • Default=N/t

schedule (dynamic [,chunk])
    • Threads grab "chunk" iterations
    • When done with iterations, thread requests next set
    • Default=1

schedule (guided [,chunk])
    • Dynamic schedule starting with large block
    • Size of the blocks shrink; no smaller than "chunk"
    • Default=1

schedule (runtime)
    • OMP_SCHEDULE
```

| Schedule Clause | When To Use |
|-----------------|---|
| STATIC | Predictable and similar work per iteration |
| DYNAMIC | Unpredictable, highly variable work per iteration |
| GUIDED | Special case of dynamic to reduce scheduling overhead |

MPI Programs

```
#include "mpi.h"
#include <stdio.h>

int main( int argc, char *argv[] )
{
    int rank, size;
    MPI_Init( &argc, &argv );
    MPI_Comm_rank( MPI_COMM_WORLD, &rank );
    MPI_Comm_size( MPI_COMM_WORLD, &size );
    printf( "I am %d of %d\n", rank, size );
    MPI_Finalize();
    return 0;
}
```

```
Process P:
send(request1,32,Q)

Process R:
send(request2, 32, Q)

Process Q:
while (true) {
    recv(received_request, 32, Any_Process);
    process received_request;
}
```

```
Process P:
send(request1, 32, Q, tag1)

Process R:
send(request2, 32, Q, tag2)

Process Q:
while (true){
    recv(received_request, 32, Any_Process, Any_Tag, Status);
    if (Status.MPI_TAG==tag1) process received_request in one way;
    if (Status.MPI_TAG==tag2) process received_request in another way;
}
```

Tags can also serve as "message type", although we should not confuse it with MPI data types

```
MPI_Init(&argc,&argv);
MPI_Comm_size(MPI_COMM_WORLD,&numprocs);
MPI_Comm_rank(MPI_COMM_WORLD,&myid);

tag=1234;source=0;destination=1;count=1;

if(myid == source){
    buffer=5678;
    MPI_Send(&buffer,count,MPI_INT,destination,tag,MPI_COMM_WORLD);
    printf("processor %d sent %d\n",myid,buffer);
}
if(myid == destination){
    MPI_Recv(&buffer,count,MPI_INT,source,tag,MPI_COMM_WORLD,&status);
    printf("processor %d got %d\n",myid,buffer);
}

MPI_Finalize();
```

MPI Basic (Blocking) Send/Receive

`MPI_SEND (start, count, datatype, dest, tag, comm)`
`MPI_RECV(start, count, datatype, source, tag, comm, status)`

Why Datatypes? ⇒ Support communication between processes on machines with different memory representations

Unsafe message passing ⇒ Depending on the availability of system buffers:

- Order the operations more carefully:

| Process 0 | Process 1 |
|-----------|-----------|
| Send(1) | Recv(0) |
| Recv(1) | Send(0) |

- Use non-blocking operations:

| Process 0 | Process 1 | Process 0 | Process 1 |
|-----------|-----------|-----------|-----------|
| Send(1) | Send(0) | Isend(1) | Irecv(0) |
| Recv(1) | Recv(0) | Irecv(1) | Send(0) |

Better:

```
if (rank == 0) {
    MPI_Isend(..., 1, tag, MPI_COMM_WORLD, &req);
    MPI_Recv(..., 1, tag, MPI_COMM_WORLD, &status);
    MPI_Wait(&req, &status);
} else if (rank == 1) {
    MPI_Recv(..., 0, tag, MPI_COMM_WORLD, &status);
    MPI_Send(..., 0, tag, MPI_COMM_WORLD);
}
```

Just six functions:

– MPI_INIT – MPI_FINALIZE – MPI_COMM_SIZE – MPI_COMM_RANK – MPI_SEND – MPI_RECV

Collective Operations in MPI:

- Collective operations are called by all processes in a communicator.
- **MPI_BCAST** distributes data from one process (the root) to all others in a communicator.
- **MPI_REDUCE** combines data from all processes in communicator and returns it to one process.
- In many numerical algorithms, **SEND/RECEIVE** can be replaced by **BCAST/REDUCE**, improving both simplicity and efficiency.

Problem with MPI: Availability

What if a process crashes?

→ If app runs 20 days or longer

 → Restart doesn't help

Checkpointing is required!

What are the obstacles of supporting heterogeneous node and fault tolerance

- Knowing N at the beginning
- Static partition of the workload

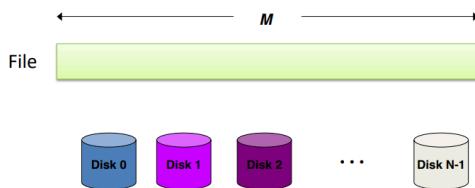
```
#include "mpi.h"
#include <math.h>
int main(int argc, char *argv[])
{
    int done = 0, n, myid, numprocs, i, rc;
    double PI25DT = 3.141592653589793238462643;
    double mypi, pi, h, sum, x, a;
    MPI_Init(&argc,&argv);
    MPI_Comm_size(MPI_COMM_WORLD, &numprocs);
    MPI_Cart_create(MPI_COMM_WORLD, 1, numprocs, 1, 1, &rc);
    h = 1.0 / (double) n;
    sum = 0.0;
    for (i = myid + 1; i <= n; i += numprocs) {
        x = h * ((double)i - 0.5);
        sum += 4.0 / (1.0 + x*x);
    }
    MPI_Reduce(&mypi, &pi, 1, MPI_DOUBLE, MPI_SUM, 0,
               MPI_COMM_WORLD);
    if (myid == 0)
        printf("pi is approximately %.16f, Error is %.16f\n",
               pi, fabs(pi - PI25DT));
}
MPI_Finalize();
return 0;
}
```

Parallel I/O Systems

Fault Tolerance

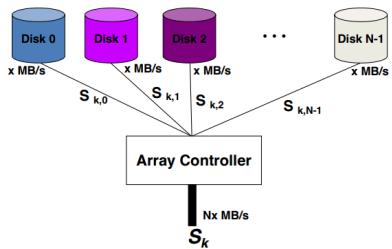
Disk Level

- Given a large file of M bytes file, and N disks
 - Find a way to distribute the file into disks.
- Objective: Read a long sequential part of the file as fast as possible



Performance Benefit

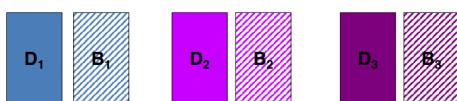
- Sequential read or write of large file*
 - application (or I/O buffer cache) reads in multiples of S bytes
 - controller performs parallel access of N disks
 - aggregate bandwidth is N times individual disk bandwidth
 - (assumes that disk is the bottleneck)



be useful

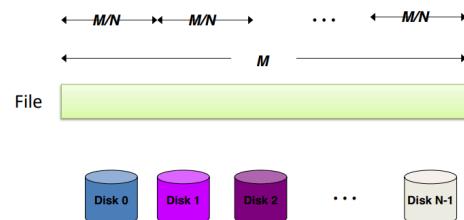
→ “Redundant Arrays of Independent Disks”

RAID Level 1



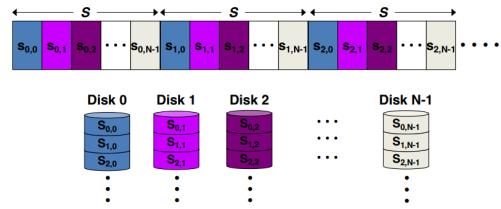
- Also known as “*mirroring*”
- To read a block:
 - read from either data disk or backup
- To write a block:
 - write both data and backup disks
- How many failure can tolerate? What is the cost for redundancy?

Solution 1: Chunk



Solution 2: Stripe

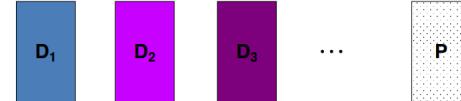
- Stripe* the data across an array of disks
 - many alternative striping strategies possible
- Example: consider a big file striped across N disks
 - stripe width* is S bytes
 - hence each *stripe unit* is S/N bytes
 - sequential read of S bytes at a time



RAID

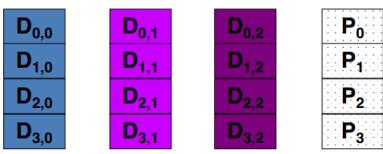
Large arrays without redundancy are too unreliable to

RAID Levels 2 & 3



- These are *bit-interleaved* schemes
- In Raid Level 2, P contains memory-style ECC
- In Raid Level 3, P contains simple parity
- Rarely used today
- Good for large data read/write, poor for small data transfer

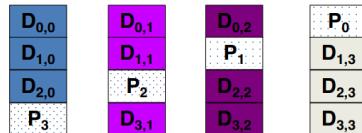
RAID Level 4



$$D_{0,0} \oplus D_{0,1} \oplus D_{0,2} = P_0$$

- *Block-interleaved parity*
- Wasted storage is small: one parity block for N data blocks
- Key problem:
 - parity disk becomes a hot spot
 - write access to parity disk on every write to any block

RAID Level 5



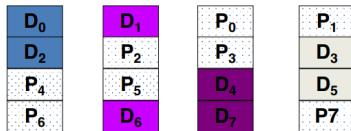
$$D_{0,0} \oplus D_{0,1} \oplus D_{0,2} = P_0$$

- *Rotated parity*
- Wastage is small: same as in Raid 4
- Parity update traffic is distributed across disks

| | RAID-0 | RAID-1 | RAID-4 | RAID-5 |
|------------------|-------------|--------------------------|-----------------------|-----------------------|
| Capacity | N | $N/2$ | $N - 1$ | $N - 1$ |
| Reliability | 0 | 1 (for sure) | 1 | 1 |
| | | $\frac{N}{2}$ (if lucky) | | |
| Throughput | | | | |
| Sequential Read | $N \cdot S$ | $(N/2) \cdot S$ | $(N - 1) \cdot S$ | $(N - 1) \cdot S$ |
| Sequential Write | $N \cdot S$ | $(N/2) \cdot S$ | $(N - 1) \cdot S$ | $(N - 1) \cdot S$ |
| Random Read | $N \cdot R$ | $N \cdot R$ | $(N - 1) \cdot R$ | $N \cdot R$ |
| Random Write | $N \cdot R$ | $(N/2) \cdot R$ | $\frac{1}{2} \cdot R$ | $\frac{N}{4} \cdot R$ |
| Latency | | | | |
| Read | D | D | D | D |
| Write | D | D | $2D$ | $2D$ |

Problem with RAID lvl 5 if failure rate not >> recovery time.

RAID Level 6



RAID provide good reliability for storage in a single node

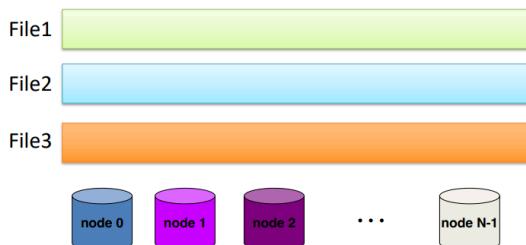
- But it **can not tolerate node/network failure**

- Tolerate 2 disk failures
- Read performance is good, while write performance is a little slower
- Parity update traffic is distributed across disks

Node Level

Problem Statement

- Given a set of large files and N machine nodes
 - Find a way to distribute the file into nodes.
- Objective:
 - Tolerate node failure
 - Read a long sequential part of the file as fast as possible
 - Write are mostly large append mode write

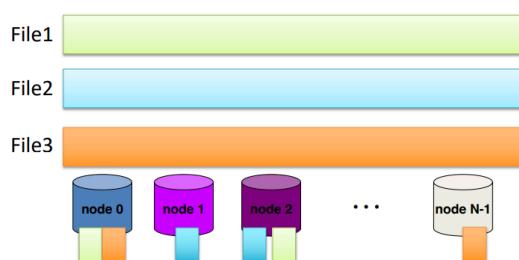


Duplicate Whole Files to Nodes

A node include all content of file.

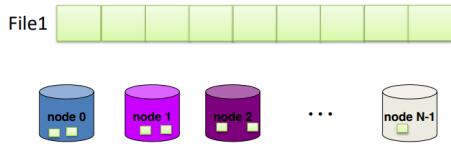
A file is duplicated in multiple nodes. (In this example:2)

Pros and cons of this strategy?



Duplicate chunks of files

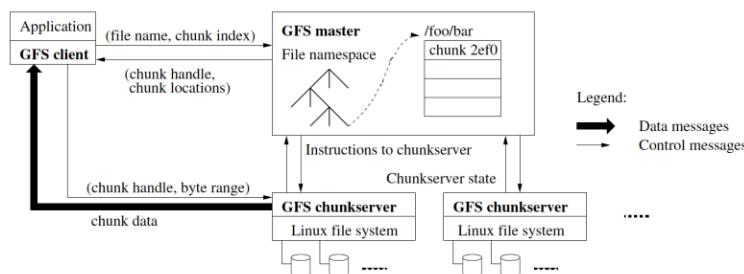
- A file is divided into fixed size chunks
- Each chunk is distributed and duplicated to a few nodes (usually 2-3)
- Pros and cons of this strategy?
 - Number of node failure tolerated?
 - Wasted space?
 - Space required to know the position of the chunk?



How to distribute 3 replicas on node/rack hierarchy?

- One replica on local node
- Second replica on a remote rack
- Third replica on same remote rack

Google File System (GFS)



Key GFS design decisions

- 3 replicas of data for fault tolerance
- Single master design for simplicity and performance
- Master does not in data paths
- Large chunk size for smaller meta data

Single master design

- Performance issue
 - Separate data path with control path, no data need to pass chunk server
 - Let's assume the system contains modest number of HUGE tiles, and the chunk size is large(e.g. 64MB)
 - Meta data for the system is not large
 - They can be put in memory of a single server, this greatly reduce the complexity of the system and could get good meta data performance.
- Single master design: Simple! With Availability and performance issues

Metadata in Memory

- The entire metadata is in main memory
- No demand paging of metadata

Data Correctness

Chunk server

A Block Server

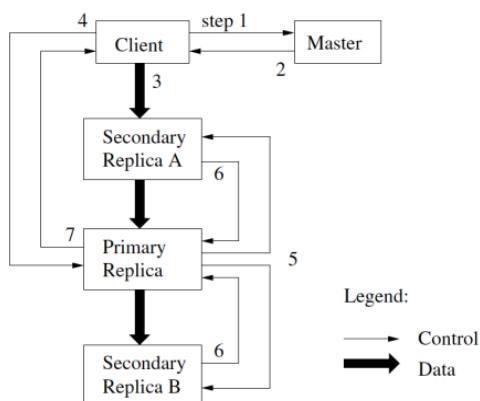
- Stores data in the local file system (e.g. ext3)
 - Each block of GFS data in a separate file in its local file system.
- Stores metadata of a block (e.g. CRC)
- Serves data and metadata to Clients

- Use Checksums to validate data
 - Use CRC32

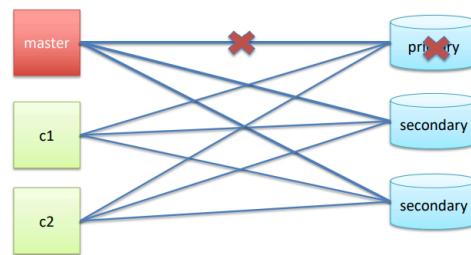
- File Creation
 - Client computes checksum per 512 bytes
 - Chunk server stores the checksum
- File access
 - Client retrieves the data and checksum from chunk server
 - If Validation fails, Client tries other replicas

GFS Write

A primary is need to control the order of the concurrent write

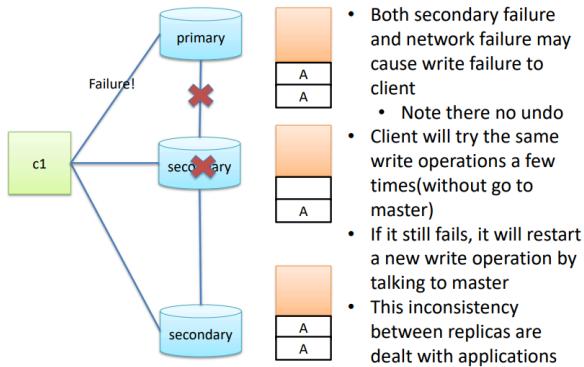


Why we need a lease?



A lease for primary is to deal with primary failure or network failure

What happens if secondary fails



Heartbeats

- Master sends heartbeat to the chunk servers
 - Once every 3 seconds
 - No heartbeat does not necessarily indicate chunk server failure
 - May be network failure

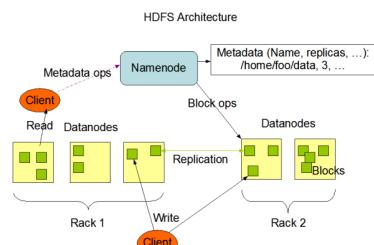
Rebalancer

- Goal: % disk full on chunk servers should be similar
 - Usually run when new chunk servers are added
 - Rebalancer is throttled to avoid network congestion

Replication Engine

- Master detects chunk server failures
 - Chooses new chunk server for new replicas
 - Balances disk usage
 - Balances communication traffic to chunk servers

HDFS Architecture



User Interface

- Commands for HDFS User:
 - hadoop dfs -mkdir /foodir
 - hadoop dfs -cat /foodir/myfile.txt
 - hadoop dfs -rm /foodir/myfile.txt
- Commands for HDFS Administrator
 - hadoop dfsadmin -report
 - hadoop dfsadmin -decommission datanodename
- Web Interface
 - http://host:port/dfshealth.jsp

Example: Count word occurrences (pseudo code)

```

map(String input_key, String input_value):
    // input_key: document name
    // input_value: document contents
    for each word w in input_value:
        EmitIntermediate(w, "1");

reduce(String output_key, Iterator intermediate_values):
    // output_key: a word
    // output_values: a list of counts
    int result = 0;
    for each v in intermediate_values:
        result += ParseInt(v);
    Emit(AsString(result));
  
```

MapReduce

- Doing a small piece of work, complete it and get the next piece
 - Input from the global file system, similar to the MPI_Bcast
 - Write a `map()` function that transfers the input piece into partial output piece
- No N is required to know at the beginning
 - Data specified by files in distributed file system
- Need a reduce support similar to MPI_Reduce
 - Write a user defined `reduce()` file

GFS design make an tradeoff

- Trade ideal consistency for system simplicity
- Clients may observe
 - Reading stale data
 - Duplicate append records
- For applications such as search engine
 - This may be tolerable
- Understanding GFS design as a tradeoff among
 - Consistency, fault-tolerance, performance, simplicity of design

Map output

Example

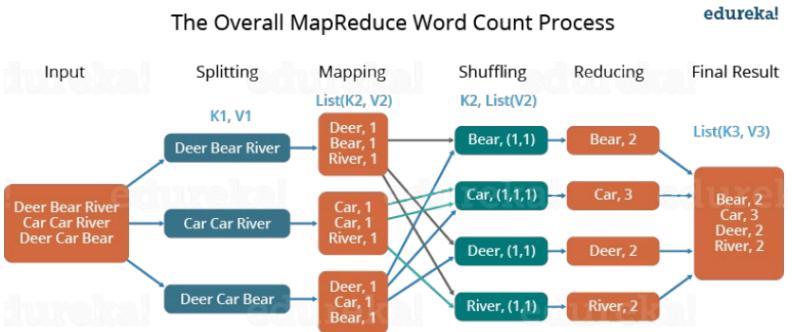
- Page 1: the weather is good
- Page 2: today is good
- Page 3: good weather is good.

Reduce Output

- Worker 1:
– (the 1)
- Worker 2:
– (is 3)
- Worker 3:
– (weather 2)
- Worker 4:
– (today 1)
- Worker 5:
– (good 4)

Parallelism

- map() functions run in parallel, creating different intermediate values from different input data sets
- reduce() functions also run in parallel, each working on a different output key
- All values are processed *independently*
- Bottleneck: reduce phase can't start until map phase is completely finished.



Why MapReduce is successful?

- Programmers only need to write serial code
 - Shared/private code in OpenMP
 - message orders in MPI
- Automatically parallel and distributed computing
- Automatic load-balancing and fault-tolerance

Limitations of MapReduce

- Only Map and Reduce
- Iterative MapReduce
 - Intermediate result on disk
 - Many I/O operations, poor performance
- Unable to support interactive query efficiently
 - No mechanism to share data in memory between MapReduce jobs

Reduce Input

- Worker 1:
– (the 1)
- Worker 2:
– (is 1), (is 1), (is 1)
- Worker 3:
– (weather 1), (weather 1)
- Worker 4:
– (today 1)
- Worker 5:
– (good 1), (good 1), (good 1), (good 1)

Transferring data with file system



- Iterative MapReduce is required when you need to "sum up" results more than 1 times
- The MapReduce model is stateless, causing performance problems for iterative MapReduce programs

Expressing an Iteration in MapReduce

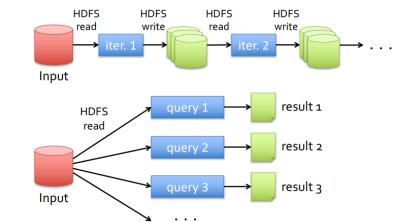
- for each point in the point set
 - find the closest mean i , assign the point to cluster i , forming new K clusters
- calculate new means for each cluster

```

map(point p):
  for each mean in K means
    find the closest mean M to this point
    emit_intermediate(index of M, p)

reduce(index of cluster i, list of points lp):
  calculate the new mean M' for points in lp
  emit(M')
  
```

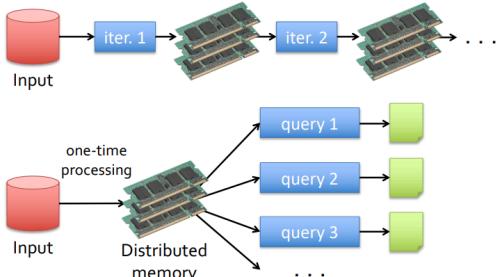
MapReduce Use File to Transfer data



The copy, serialization and disk I/O makes it slow

In Memory Computing

Use memory to store data



10-100 times faster than the disk approach

Issues

- Is memory large enough to contain all the data?
- How about the cost comparing to disks?
- Fault-tolerance with memory?
- How to represent the data in memory efficiently?

Two traditional abstractions

- Distributed Shared Memory
 - Extend the virtual memory
read/write #address
- Distributed KV stores
 - Value = Get(key)
 - Put(key, value)

How to tolerate fault in distributed k-v stores?

A Log Sample

| TID | T1 | T1 | T1 | T2 | T2 | T2 | T3 |
|-----|-------|------|--------|-------|------|--------|-------|
| OLD | A=0 | B=0 | | A=100 | B=50 | | A=80 |
| NEW | A=100 | B=50 | COMMIT | A=80 | B=70 | COMMIT | A=110 |
| | | | | | | | |

```
begin // T1
A = 100
B = 50
commit // A=100; B=50

begin // T2
A = A - 20
B = B + 20
commit // A=80; B=70

begin // T3
A = A + 30 // A=110
--CRASH--
```

Fault tolerance mechanism of DSM or Key-value stores

- Replica
 - High cost
 - Big time overhead
- Log
 - Big time overhead
- Remote memory access
 - 10-100 times slower than local memory accesses
 - 1-10us vs 100ns

A more lightweight fault tolerance mechanism for in-memory computing is desired

Five Operations Involved

- **1. begin:** allocate a new transaction ID
- **2. write variable:** append an entry to the log
- **3. read variable:** scan the log looking for last committed value
 - As an aside: how to see your own updates?
 - Read uncommitted values from your own TID

Five Operations Involved

Performance Problem of Log-only Approach

- **Write performance** is probably good
 - Sequential writes, instead of random
 - Need to write to disk **twice** instead of **once**
- **Read performance** is **terrible!**
 - Scan the log for every read!
- **Recovery** is instantaneous
 - Nothing to do

- **4. Commit:** write a commit record
 - Expectedly, writing a commit record is the "commit point" for action, because of the way read works (looks for commit record)
 - However, writing log records better be **all-or-nothing**
 - One approach, from last time: make each record fit within one sector
- **5. Abort:** Do nothing or write a record?
 - Could write an abort record, but not strictly needed
 - *Recover from a crash:* do nothing

RDD (Resilient Distributed Datasets)

Solution

RDD (Resilient Distributed Datasets)

- Based on data sets, instead of single data
- Generated with deterministic coarse grained operations (map, filter, join etc.)
- Immutable
- To modify data, create new datasets by transformation

Efficient Fault Tolerance

- Once the data is generated deterministically, and immutable
 - Recover the data from “recompute”
 - Only need to log the sequences to generate rdd. small cost if there's no fault happens

```
messages = textFile(...).filter(_.contains("error"))
           .map(_.split("\t")(2))
```



How to create RDDs

- From Collection


```
val a = sc.parallelize(1 to 9, 3)
```
- from File


```
val b = sc.textFile("/path/to/file")
```
- from other RDDs


```
val c = b.map(line => line.split(","))
```

Spark Operators

| | | |
|--|---|---|
| Transformations (define a new RDD) | map filter sample groupByKey reduceByKey sortByKey | flatMap union join cogroup cross mapValues |
| Actions (return a result to driver program) | collect reduce count save lookupKey | |

Spark Programming

Scala - Functions

- def helloWord()={ println("Hello World")}
- def addOne(x:Int) = {


```
x+1
```

 }
- def addOne(x:Int) = x+1
- val addOne = (x:Int) => x + 1

RDD operations

- Transformation: generate new RDDs from existing RDDs
 - map, filter, groupBy, sort, distinct, sample...
- Action: return a value from RDD
 - **reduce**, count, collect, first, foreach..., ,

Map

- Map each element in source RDD to ONE element in destination RDD

```
scala> val a = sc.parallelize(1 to 9, 3)
scala> val b = a.map(x => x*2)
scala> a.collect
res10: Array[Int] = Array(1, 2, 3, 4, 5, 6, 7, 8, 9)
scala> b.collect
res11: Array[Int] = Array(2, 4, 6, 8, 10, 12, 14, 16, 18)
reduceByKey
```

- Reduce for key-value list
- Values with the same keys are reduced, and form a new key-value

```
scala> val a = sc.parallelize(List((1,2),(3,4),(3,6)))
scala> a.reduceByKey((x,y) => x + y).collect
res7: Array[(Int, Int)] = Array((1,2), (3,10))
```

Use of _

- val a = (1 to 9)
- Val b = a.map(x => x*2)

is the same as val b = a.map(_ * 2)
- val a = sc.parallelize(List((1,2),(3,4),(3,6))) **a.reduceByKey((x,y) => x + y).collect** is the same as **a.reduceByKey(_ + _).collect**

flatMap

- One element in source rdd will be transformed to multiple elements(0..n) in destination rdd

```
scala> val a = sc.parallelize(1 to 4, 2)
scala> val b = a.flatMap(x => 1 to x)
scala> b.collect
res0: Array[Int] = Array(1, 1, 2, 1, 2, 3, 1, 2, 3, 4)
```

Join

(k,v1) join (k, v2) → (k, (v1, v2))

```
scala> val kv1=sc.parallelize(List(("A",1),("B",2),("C",3)))
scala> val kv3=sc.parallelize(List(("A",4),("A",5),("B",3),("D",30)))
scala> kv1.join(kv3).collect
Array[(String, (Int, Int))] = Array((A,(1,4)), (A,(1,5)), (B,(2,3)))
scala> kv1.union(kv3).collect
Array[(String, (Int, Int))] = Array((A,1), (B,2), (C,3), (A,4), (A,5), (B,3), (D,30))
```

```
scala> val c = sc.parallelize(1 to 10)
scala> c.reduce((x, y) => x + y)
```

res4: Int = 55

Output

```
• saveAsTextFile/saveAsSequenceFile/saveAsObjectFile
```

```
scala> val a = sc.parallelize(List((1,2),(3,4),(3,6)))
scala> a.reduceByKey((x,y) => x + y).saveAsTextFile("test")
```

Join

(k,v1) join (k, v2) → (k, (v1, v2))

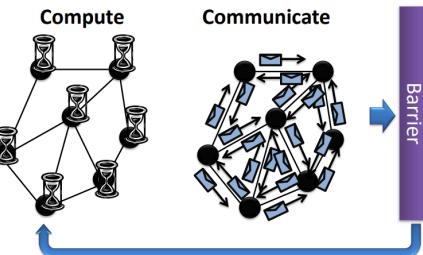
```
scala> val kv1=sc.parallelize(List(("A",1),("B",2),("C",3)))
scala> val kv3=sc.parallelize(List(("A",4),("A",5),("B",3),("D",30)))
scala> kv1.join(kv3).collect
Array[(String, (Int, Int))] = Array((A,(1,4)), (A,(1,5)), (B,(2,3)))
scala> kv1.union(kv3).collect
Array((A,1), (B,2), (C,3), (A,4), (A,5), (B,3), (D,30))
```


Graph Analysis

Spark isn't suited for extreme-scale graph processing

- Slow and memory inefficient

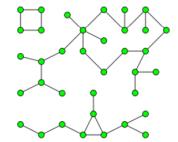
The BSP model does not support asynchronous computing



Weakly Connected Components

Algorithm

- Give each vertex an ID and look at all edges and find the vertex ID who holds the highest value.
- Send the greatest vertex ID to all the edges.
- Process incoming messages and compare it against the value you hold.
- Repeat the process until there are no messages being sent.



What happens if this algorithm is executed asynchronously?

Asynchronous processing

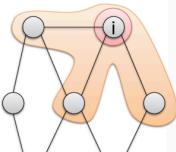
The GraphLab Abstraction

Vertex-Programs directly **read** the neighbors state

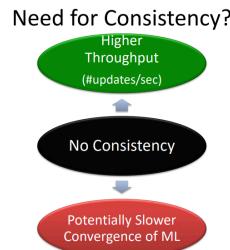
```
GraphLab_PageRank(i)
    // Compute sum over neighbors
    total = 0
    foreach( j in in_neighbors(i)):
        total = total + R[j] * wji

    // Update the PageRank
    R[i] = 0.15 + total

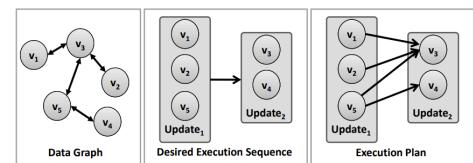
    // Trigger neighbors to run again
    if R[i] not converged then
        foreach( j in out_neighbors(i)):
            signal vertex-program on j
```



Ensuring Race-Free Code in Async processing:

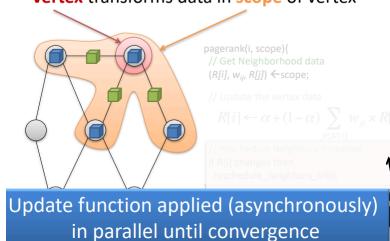


An execution plan example



Update Functions

User-defined program: applied to **vertex** transforms data in **scope** of vertex



Random Partitioning

- Both GraphLab 1 and Pregel proposed Random (hashed) partitioning for Natural Graphs

For p Machines:

$$\mathbb{E} \left[\frac{|Edges Cut|}{|E|} \right] = 1 - \frac{1}{p}$$

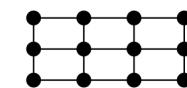
10 Machines \rightarrow 90% of edges cut
100 Machines \rightarrow 99% of edges cut!

Graph partitioning on power law graphs and GAS model

Needed because Graphs can become unhandable huge.

Assumptions of Graph-Parallel Abstractions

Idealized Structure



- Small neighborhoods**
 - Low degree vertices
- Vertices have similar degree
- Easy to partition

Natural Graph



- Large Neighborhoods**
 - High degree vertices
- Power-Law** degree distribution
- Difficult to partition**

In Summary

GraphLab 1 and Pregel are not well suited for natural graphs

Problem:

High Degree Vertices Limit Parallelism



Edge information too large for single machine



Touches a large fraction of graph (GraphLab 1)



Produces many messages (Pregel)



Sequential Vertex-Updates



Asynchronous consistency requires heavy locking (GraphLab 1)



Synchronous consistency is prone to stragglers (Pregel)

- Poor performance on high-degree vertices
- Low Quality Partitioning



Minimizing Communication in GraphLab 2: Vertex Cuts

- Distribute a single vertex-update
 - Move computation to data
 - Parallelize high-degree vertices
- Vertex Partitioning
 - Simple online approach, effectively partitions large power-law graphs

A Common Pattern for Vertex-Programs

```
GraphLab_PageRank(i)
// Compute sum over neighbors
total = 0
foreach( j in in_neighbors(i)):
    total = total + R[j] * wji

// Update the PageRank
R[i] = 0.1 + total

// Trigger neighbors to run again
if R[i] not converged then
    foreach( j in out_neighbors(i))
        signal vertex-program on j
```

Gather Information About Neighborhood

Update Vertex

Signal Neighbors & Modify Edge Data

PageRank in PowerGraph

$$R[i] = 0.15 + \sum_{j \in \text{Nbrs}(i)} w_{ji} R[j]$$

PowerGraph_PageRank(i)

```
Gather(j → i) : return wji * R[j]
sum(a, b) : return a + b;
```

```
Apply(i, Σ) : R[i] = 0.15 + Σ
```

Scatter(i → j) :

if $R[i]$ changed then trigger j to be recomputed

PowerLyra

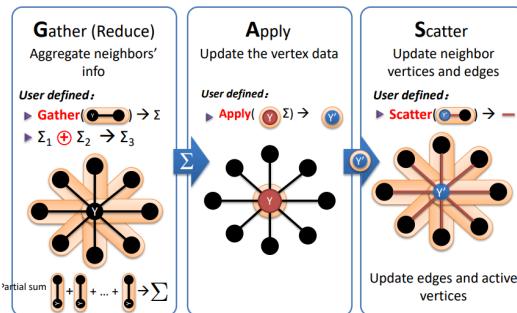
- PowerGraph partitions graph for high degree vertices, but is not friendly to the low degree vertices(locality)
- Optimize the low degree vertices' locality
 - Avoids unnecessary communications

A vertex-cut minimizes
machines per vertex

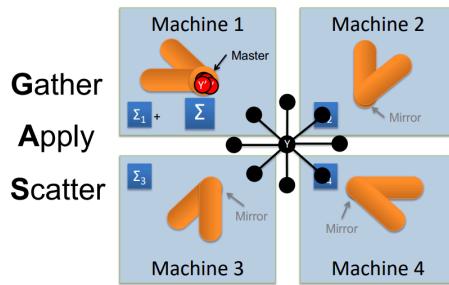
Percolation theory suggests Power Law graphs can be split by removing only a small set of vertices [Albert et al. 2000]

Small vertex cuts possible!

GAS



Distributed GAS



Constructing Vertex-Cuts

- Goal:** Parallel graph partitioning on ingress
- GraphLab 2 provides three simple approaches:
 - Random Edge Placement**
 - Edges are placed randomly by each machine
 - Good theoretical guarantees
 - Greedy Edge Placement with Coordination**
 - Edges are placed using a shared objective
 - Better theoretical guarantees
 - Oblivious-Greedy Edge Placement**
 - Edges are placed using a local objective

Out of core graph processing

The key of out-of-core graph systems is the sequential access to graph data

Generalized Graph Processing

```
for each vertex v
  if v is active
    for each edge e to v
      gather update along e
      apply updated value to v
      Vertex 1.0 1.0 1.0 1.0
```

| | | | | | | |
|------|-------|-------|-------|-------|-------|-------|
| Edge | 1 → 2 | 1 → 3 | 2 → 1 | 2 → 3 | 3 → 2 | 4 → 2 |
| | 2 → 1 | 2 → 4 | 1 → 2 | 1 → 3 | 4 → 3 | 2 → 4 |
| | 3 → 2 | 4 → 3 | 2 → 3 | 2 → 4 | | |
| | 4 → 2 | 4 → 3 | | | | |

Can we use just one copy of the edge list?
Do we need all vertices in memory?

Out-of-core Graph Processing

- Processing large graphs with limited memory
- Low cost graph analysis engine

Possible Solutions

1. Use SSD as a memory extension? (Hadoop, MapReduce)
2. Compress the graph structure to fit into RAM? (WebGraph framework)
Too many small updates per second / sec.
3. Cluster the graph and handle each cluster separately in RAM?
Expensive. The number of inter-cluster edges is big.
4. Caching of hot nodes?
Unpredictable performance.

Generalized Graph Processing

```
for each vertex v
  if v is active
    for each edge e to v
      gather update along e
      apply updated value to v
```

| | | | |
|------|-------|-------|-------|
| Edge | 1 → 2 | 1 → 3 | 2 → 1 |
| | 2 → 1 | 2 → 4 | 1 → 2 |
| | 3 → 2 | 4 → 3 | 3 → 2 |
| | 4 → 2 | 4 → 3 | 4 → 2 |

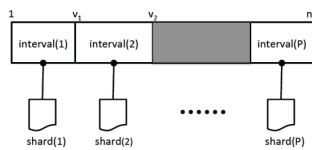
Gather&Apply can be separated from Scatter

GRAPHCHI – Parallel Sliding Window

PSW: Shards and Intervals

| |
|------------|
| 1. Load |
| 2. Compute |
| 3. Write |

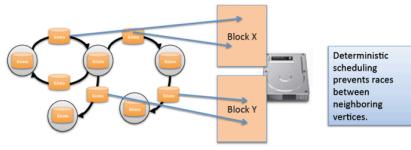
- Vertices are numbered from 1 to n
 - P intervals, each associated with a shard on disk.
 - sub-graph = interval of vertices



PSW: Execute updates

| |
|------------|
| 1. Load |
| 2. Compute |
| 3. Write |

- Update-function is executed on interval's vertices
- Edges have pointers to the loaded data blocks
- Changes take effect immediately → asynchronous.



Graphchi processing model

for each vertex partition (Intervals)
 load all vertex and in/out edges
 for each vertex v
 if v is active
 for each edge e to v
 gather update along e
 apply updated value to v
 for each vertex v
 if v has update
 for each edge e from v
 scatter update along e by attaching the info
 Write all in/out edges

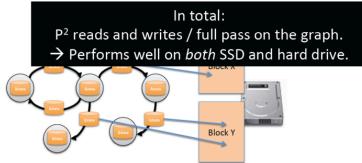
Dual sliding window to allow sequential access to edges

Remove the need to sort in preprocessing

2-D partition to allow selective scheduling on edges; In place apply to avoid shuffle

PSW: Commit to Disk

- In write phase, the blocks are written back to disk
 - Next load-phase sees the preceding writes → asynchronous.



GraphChi is good

- But have two major drawbacks

- P^2 I/O are required
 - Poor performance when P is large
- High preprocessing time because sorting is required for each shard

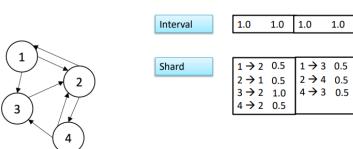
Example



How is gather and scatter implemented?

Example with (simplified)

PageRank

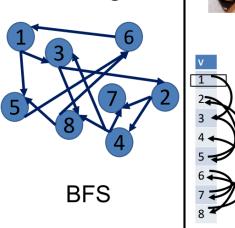


XSTREAM – Edge Centric

X-Stream : Edge-Centric model

| | |
|--|---------|
| for each vertex v if v has update for each edge e from v scatter update along e | Scatter |
| for each edge e If e.src has update scatter update along e | Scatter |

Edge Centric Processing



| SOURCE | DEST |
|--------|------|
| 1 | 3 |
| 1 | 5 |
| 2 | 7 |
| 2 | 4 |
| 3 | 2 |
| 3 | 8 |
| 4 | 3 |
| 4 | 7 |
| 4 | 8 |
| 5 | 6 |
| 5 | 1 |
| 6 | 1 |
| 6 | 8 |
| 7 | 3 |
| 7 | 4 |
| 7 | 6 |
| 8 | 5 |
| 8 | 6 |

Vertex-centric Scatter-Gather: Random Access Bandwidth

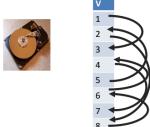
- Edge-centric Scatter-Gather
 - #of scatters * Edge data / sequential access bandwidth
- Sequential Access Bandwidth >> Random Access Bandwidth
- Few scatter-gather iterations for real world graphs
 - Give an example that is not suitable for this system?

Vertex-Centric → Edge-Centric

Streaming Partitions

Streaming Partitions

- Problem: still have random access to vertex set



- Solution: partition the graph into streaming partitions

Streaming Partitions

- A streaming partition is
 - A subset of the vertices that fits in RAM
 - All edges whose source vertex is in that subset
 - No requirement on quality of the partition

The issues of X-Stream

- X-Stream
 - Separate gather/scatter phase, with shuffle between them, introduce many extra I/O operations
 - Poor performance for BFS and WCC-like algorithms
 - Only a small portion of edges are touched for some iterations, but they would scan all edges

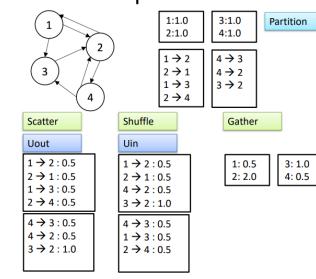
Partition Edges and vertices

| SOURCE | DEST |
|--------|------|
| 1 | 5 |
| 4 | 7 |
| 2 | 7 |
| 4 | 3 |
| 3 | 8 |
| 2 | 4 |
| 1 | 3 |
| 3 | 2 |

| V1 | V2 |
|----|----|
| 1 | 5 |
| 2 | 7 |
| 3 | 8 |
| 4 | 3 |
| 5 | 6 |
| 6 | 8 |
| 7 | 5 |
| 8 | 1 |

Edges in disks/SSDs

Example on X-Stream



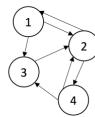
GRIDGRAPH – 2D Partition

GridGraph

- A 2-level partitioning graph data structure
 - Improve locality
 - Enable more effective streaming on both edges and vertices
- A Streaming-Apply model
 - Combines the gather-apply-scatter operations
- A programming abstraction
 - To enable the streaming-apply model

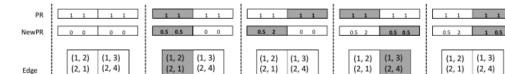
The Grid Representation

- 2-level hierarchical partitioning
 - 1st dimension
 - partition the vertices into chunks(1,2)(3,4)
 - partition the edges into shards by source vertex((1,2),(2,1),(1,3),(2,4))
 - 2nd dimension
 - partition the shards into blocks by destination vertex



| | |
|--------|--------|
| (1, 2) | (1, 3) |
| (2, 1) | (2, 4) |
| (3, 2) | (4, 3) |
| (4, 2) | |

The Streaming-Apply Model



Dual Sliding Windows:

- 1 read window + 1 write window
- Contents of active windows can be fit into RAM

Each iteration, just read edge list once, has opportunity to skip some locks

```
Algorithm 1 Edge Block Streaming
for j ← 1, P do
  for i ← 1, P do
    if ChunksActive(i) then
      StreamEdgeBlock(i,j);
    end if
  end for
end for
```

Streaming-Apply Processing Model

- Vertex accesses are aggregated
 - Good locality
 - Only 2 partitions of vertices are accessed within each edge block
- On-the-fly updates onto vertices
 - Reduces I/O
 - Enables asynchronous implementation of algorithms (like WCC, etc.) which converges faster

GridGraph processing model

```
for each vertex partition
  load all vertex and in/out edges

  for each edge e
    gather update along e
    read e.src
    apply updated value to e.dst

Write the unnecessary vertex partition back to disk
```

Programming Abstraction

Algorithm 2 Edge Streaming Interface

```
function STREAMEDGES(Fe, F)
  Sum = 0
  for all block do
    for all edge ∈ block do
      if F(edge.source) then
        Sum += Fe(edge)
      end if
    end for
  end for
  return Sum
end function
```



StreamEdges(Fe, F):
Fe: process the edge, which is allowed to side-effect related vertex data
F: check whether the edge should be processed (by looking at the source vertex)

Algorithm 3 Vertex Streaming Interface

```
function STREAMVERTICES(Fv, F)
  Sum = 0
  for all vertex do
    if F(vertex) then
      Sum += Fv(vertex)
    end if
  end for
  return Sum
end function
```

StreamVertices(Fv, F):
Fv: process the vertex
F: check if the vertex should be processed

BFS

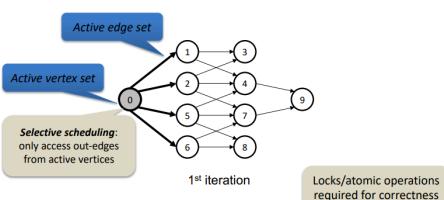
```
Algorithm 4 BFS
function ISACTIVE(v)
  return ActiveIn[v]
end function
function VISIT(e)
  if CAS(&Parent[e.dest], -1, e.source) then
    ActiveOut[e.dest] = 1
    return 1
  end if
  return 0
end function
ActiveVertices = 1
Parent = {-1, ..., -1}
ActiveIn = {0, ..., 0}
Parent[start] = start
ActiveIn[start] = 1
while ActiveVertices > 0 do
  ActiveOut = {0, ..., 0}
  ActiveVertices = StreamEdges(VISIT, ISACTIVE)
  Swap(ActiveIn, ActiveOut)
end while
```

PageRank Implementation

```
Algorithm 3 PageRank
function CONTRIBUTE(e)
  Accum(&NewPR[e.dest], PR[e.source])
end function
function COMPUTE(v)
  NewPR[v] = 1 - d + d × NewPR[v]
  return |NewPR[v] - PR[v]|
end function
d = 0.85
PR = {1, ..., 1}
Converged = 0
while ¬Converged do
  NewPR = {0, ..., 0}
  StreamEdges(Contribute)
  Diff = StreamVertices(Compute)
  Swap(PR, NewPR)
  Converged = Diff / V ≤ Threshold
end while
```

Top Down vs. Bottom Up

BFS Example (top-down model)



function top-down-step(vertices, frontier, next, parents)

```
for v ∈ frontier do
  for n ∈ neighbors[v] do
    if parents[n] = -1 then
      parents[n] ← v
      next ← next ∪ {n}
    end if
  end for
end for
```

Heuristics for mode switch

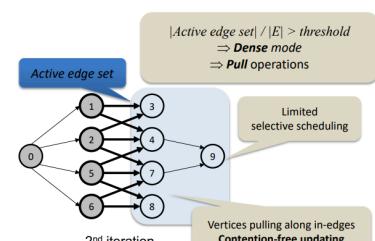
Dual mode updates proposed in shared-memory systems (LigrapPPUPP-13)

```
|Active edge set| / |E| < threshold
⇒ Sparse mode
⇒ Push operations
```

```
|Active edge set| / |E| > threshold
⇒ Dense mode
⇒ Pull operations
```

- Top-down
 - Less active edges – sparse – push
- Bottom-up
 - More active edges – dense – pull

BFS(Bottom-up mode)



Bottom-up BFS

function bottom-up-step(vertices, frontier, next, parents)

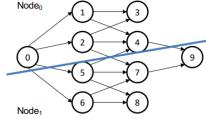
```
for v ∈ vertices do
  if parents[v] = -1 then
    for n ∈ neighbors[v] do
      parents[v] ← n
      next ← next ∪ {v}
    end for
  end if
end for
```

GEMINI – Distributed push/pull

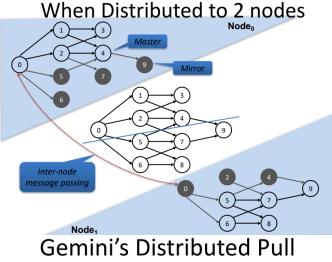
We Propose: *Gemini*

- Build **scalability** on top of **efficiency**
 - Avoid unnecessary “distributed” side-effects
 - Optimize computation on partitioned graphs
- Shift of design focus
 - Designed for distributed, but **computation-centric**
 - Modern clusters have fast interconnects
 - Computation-communication overlap in place
 - Efficiency optimizations
 - Adaptive push-/pull-style computation
 - Hierarchical chunk-based partitioning
 - Scalability optimizations
 - Locality-aware chunking
 - Chunk-based work-stealing

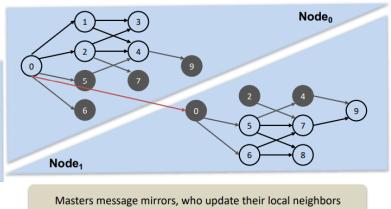
Distributed Dual-mode Computation



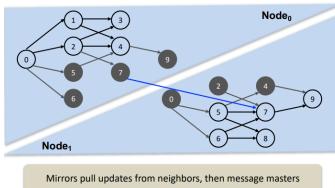
When Distributed to 2 nodes



Gemini’s Distributed Push

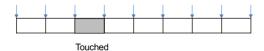


Gemini’s Distributed Pull



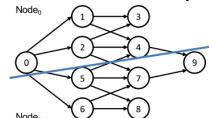
More Benefits of Chunking

- No need to convert vertex IDs (global \Leftrightarrow local)
 - Simple bookkeeping of partition information
 - $O(p)$ chunk boundaries
 - Simple management of vertex data arrays
 - Just allocate entire array in shared memory



- Can be applied recursively at different levels

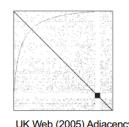
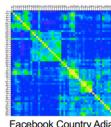
Gemini’s Choice for Graph Partition



- **Chunking**
 - Divide vertex set V into p contiguous chunks
 -
 - Dual-mode edge data distributed accordingly

Why Chunk-based Partitioning?

- It preserves locality!
 - Fact: **locality** exists in many real-world graphs
 - E.g., WebGraph[WWW’04], BLP[WSOM’13]
 - Vertices “semantically” ordered
 - Small distance between source and destination vertices
- Preprocessing affordable when vertices unordered
 - E.g., BFS[Algorithms’09], LLPI[WWW’11]



Locality-aware Chunking

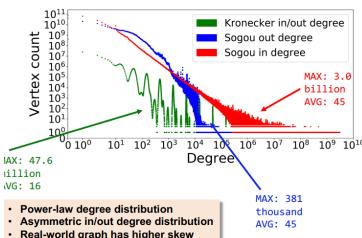
- Where to draw partition lines?
 -
 -
 - Balancing by edges?

Chunk size affects random access efficiency!

- Gemini considers both vertex and edge
 - Edge: the amount of work to be processed
 - Vertex: the processing speed of work (locality)
 - Hybrid metric: $\alpha \cdot |V_i| + |E_i|$
 - α set to 8/(p-1) currently default value

Shentu - Designed for Extreme Scale Graph

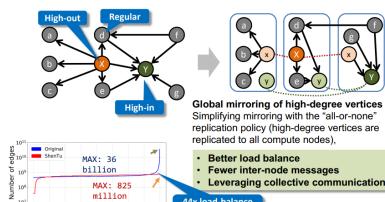
Problem 1: Severe Load Imbalance



Problem 3: Heterogeneous Architecture

- MPE (Management Processing Element)
 - General-purpose processor
 - Low memory bandwidth
- CPE (Compute Processing Element)
 - High memory bandwidth
 - Limited functionality
- Shared memory between MPE and CPE

Solution 1: Degree Aware Messaging



Problem 2a: Too Many Small Messages

- Traditional distributed graph computing model (1D)
 - $O(N \times N)$ messages
 - Not suitable to extreme-scale processing
 - Too expensive on 40,000+ nodes
 - # nodes >> average degree: little opportunity for combining messages
- Hierarchical interconnect
 -

Solution 2a: Supernode Routing

- Stage 1
 - Messages grouped by destination supernodes
 - Combined messages sent to relay node in destination supernode
- Stage 2
 - Messages grouped by destination nodes within same supernode
 - Combined messages sent to final destination node

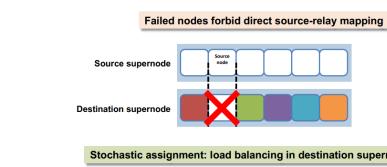
Solution 3b: CPE Cluster On-chip Sorting



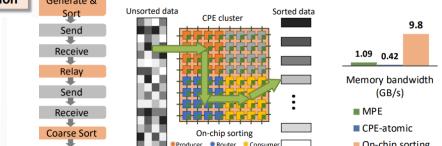
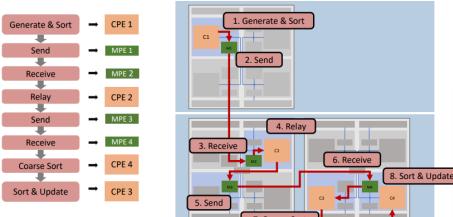
Problem 2a: Too Many Small Messages

- Traditional distributed graph computing model (1D)
 - $O(N \times N)$ messages
 - Not suitable to extreme-scale processing
 - Too expensive on 40,000+ nodes
 - # nodes >> average degree: little opportunity for combining messages
- Hierarchical interconnect
 -

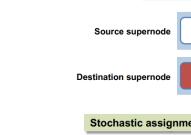
Solution 2b: Topology-aware Supernode Relaying



Solution 3a: Heterogeneous Task Mapping



Solution 3b: CPE Cluster On-chip Sorting



Stream Processing – STORM

Core Concepts and Architecture

1. **Stream:**
 - o The core abstraction in Storm is the "stream," which is an unbounded sequence of tuples (data records).
 - o Streams are the primary input/output data format for all operations in a Storm topology.
2. **Tuple:**
 - o A tuple is a named list of values, essentially the smallest unit of data processed in a Storm application.
 - o Tuples can contain any number of fields, such as integers, strings, or objects.
3. **Topology:**
 - o A topology is a directed acyclic graph (DAG) that represents the flow of data between computations in Storm.
 - o A topology runs indefinitely until terminated, processing streams of data in real-time.
4. **Spouts:**
 - o Spouts are the sources of streams in a Storm topology.
 - o They emit data tuples into the topology, typically pulling data from external sources like message queues (e.g., Kafka), databases, or APIs.
5. **Bolts:**
 - o Bolts process tuples from spouts or other bolts. They are the building blocks for processing logic.
 - o Bolts can perform various tasks such as filtering, aggregation, joining streams, or interacting with databases.
6. **Stream Grouping:**
 - o Determines how tuples are routed between bolts.
 - o Examples:
 - **Shuffle Grouping:** Distributes tuples randomly.
 - **Fields Grouping:** Routes tuples to bolts based on specific fields.
 - **All Grouping:** Sends each tuple to all bolts.
 - **Global Grouping:** Sends all tuples to a single bolt.

Key Features

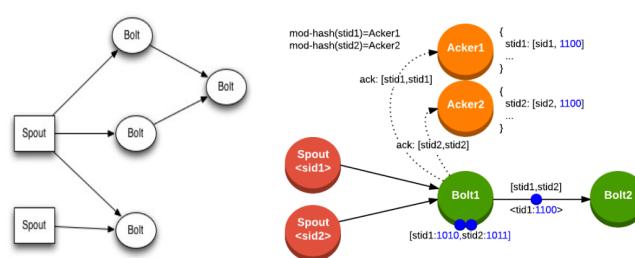
1. **Fault Tolerance:**
 - o Storm ensures reliability through an acknowledgment system. If a tuple fails to process, it can be replayed.
 - o Worker nodes can restart automatically in case of failure.
2. **Scalability:**
 - o Storm can handle a high throughput of messages by parallelizing computation.
 - o Parallelism is achieved by specifying the number of executors and tasks for spouts and bolts.
3. **Low Latency:**
 - o Designed for sub-second processing times, Storm is ideal for applications requiring immediate insights.
4. **Language Agnostic:**
 - o Storm's components can be written in any programming language that supports JSON serialization.

Stream Grouping

- Shuffle Grouping: Randomly
- Fields Grouping: Choose according to tuple field's value
- All Grouping: all tasks
- Global Grouping: send to tasks with the minimum id

Fault tolerance in Stream Processing

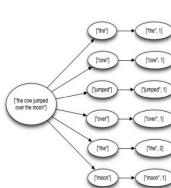
- Per element processing(at least once) - Storm
- Mini-batch(Exactly once) – SparkStream
- Per element processing(Exactly once) - Flink



1 Define the spout

Complete processing of a message

- When spouts emit tuples, specify a unique message ID.
- When bolts are done processing the input tuple, ack or fail the input tuple.



TestWordSpout

```
public void nextTuple() {
    Utils.sleep(100);
    final String[] words = new String[] {"nathan", "mike", "jackson", "golda", "bertels"};
    final Random rand = new Random();
    final String word = words[rand.nextInt(words.length)];
    _collector.emit(new Values(word));
}
```

emit() is to push the data to downstream bolts

WordCount in Storm

```
public static class WordCount implements IBasicBolt {
    Map<String, Integer> counts = new HashMap<String, Integer>();

    public void prepare(Map conf, TopologyContext context) {
    }

    public void execute(Tuple tuple, BasicOutputCollector collector) {
        String word = tuple.getString(0);
        Integer count = counts.get(word);
        if(count==null) count = 0;
        count++;
        counts.put(word, count);
        collector.emit(new Values(word, count));
    }

    public void cleanup() {
    }
}
```

Running Storm

```
LocalCluster cluster = new LocalCluster();

Map conf = new HashMap();
conf.put(Config.TOPOLOGY_DEBUG, true);

cluster.submitTopology("demo", conf, builder.createTopology());

cluster.shutdown();
```

How to track if the message has been processed?

- We can follow the message flow to see if there are outstanding messages
- Use a data structure to record the generations/acknowledgement of messages
 - Do we need a new entry for each generated message?



Maintain the pending message

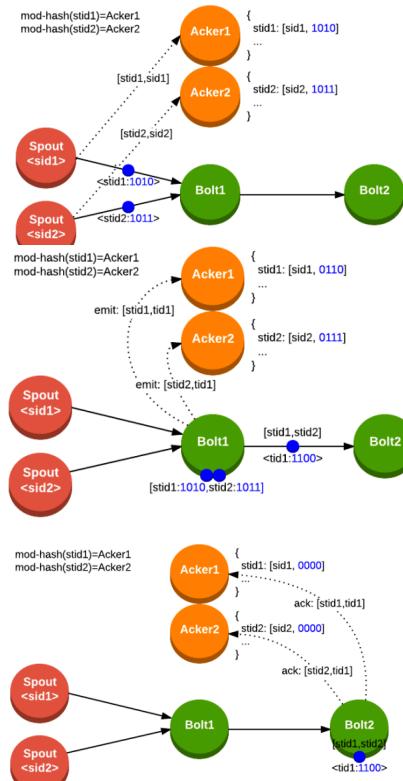
- Message 1 , acknowledged
- Message 2, pending
- ...

Process and acknowledge message 1

Process and acknowledge message 2

High memory overhead, also does not know who is the root message of them

A better way is to have an entry for each source message



3. Define the topology

ExclamationTopology

```
TopologyBuilder builder = new TopologyBuilder();
builder.setSpout("words", new TestWordSpout(), 10);
builder.setBolt("exclaim1", new ExclamationBolt(), 1).shuffleGrouping("words");
builder.setBolt("exclaim2", new ExclamationBolt(), 2).shuffleGrouping("exclaim1");
```

2. Define the bolts

ExclamationBolt

```
public void execute(Tuple tuple) {
    _collector.emit(tuple, new Values(tuple.getString(0) + "!!!"));
    _collector.ack(tuple);
}

public void declareOutputFields(OutputFieldsDeclarer declarer) {
    declarer.declare(new Fields("word"));
}
```

Notice the ack() here, why is it necessary?

SparkStream

Review: How Spark deal with Fault Tolerance

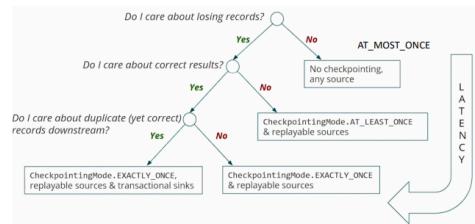
- Computation are implemented as deterministic transformations on immutable data sets

```
messages = textFile(...).filter(_.contains("error"))
           .map(_.split("\t")(2))
```



Different fault tolerance levels

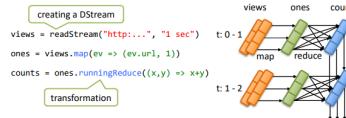
At most once, At least once, Exactly once



Example: Counting page views

Discretized Stream (DStream) is a sequence of immutable, partitioned datasets

- Can be created from live data streams or by applying bulk, parallel **transformations** on other DStreams



SparkStream WordCount code example*

```
// Create a ReceiverInputDStream on target ip:port and count the
// words in input stream of \n delimited test (e.g. generated by 'nc')
val lines = ssc.socketTextStream(args(0), args(1).toInt)
val words = lines.flatMap(_.split(" "))
val wordDStream = words.map(x => (x, 1))

// Update the cumulative count using mapWithState
// This will give a DStream made of state (which is the cumulative count of the words)
val mappingFunc = (word: String, one: Option[Int], state: State[Int]) => {
  val sum = one.getOrElse(0) + state.getOption.getOrDefault(0)
  val output = (word, sum)
  state.update(sum)
}

// Run a streaming computation as a series of
// small, deterministic batch jobs
```

Discretized Stream Processing

Run a streaming computation as a series of small, deterministic batch jobs

Store intermediate state data in cluster memory

Try to make batch sizes as small as possible to get second-scale latencies

Fault tolerance

Fault Tolerance

- Storm guarantees a message is processed (at least) once
- SparkStream uses mini-batch to provide exactly once
- Flink (and storm trident) can do per-message exactly once with checkpoint, with replayable source(kafka) and transactional sink

Comparison of different Streaming Models

Comparison of different systems

- Programming model**
 - Storm and Flink requires save and processing more
 - SparkStream is an extension of Spark
 - But Flink provide SQL like interface which is very user friendly

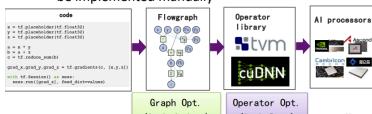
Flow control

- Push or pull**
 - Spark uses push
 - Storm and Flink uses pull, buffer overflow only at the first stage
 - Pull is also known as back pressure

Machine Learning

Machine Learning Systems

- Compile the program to flowgraph and optimize it
- Support libraries for new hardware
 - Thousands of operators which are not affordable to be implemented manually



Matrix Multiplication Loop Orders

ijk (& jik):

- 2 loads, 0 stores
- misses/iter = 1.25

```
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    for (k=0; k<n; k++) {
      sum = 0.0;
      for (l=0; l<n; l++) {
        sum += a[i][j] * b[k][l];
        c[i][l] = sum;
      }
    }
  }
}
```

kij (& ikj):

- 2 loads, 1 store
- misses/iter = 0.5

```
for (i=0; i<n; i++) {
  for (j=0; j<n; j++) {
    for (k=0; k<n; k++) {
      r = a[i][k];
      for (l=0; l<n; l++) {
        sum += r * b[k][l];
      }
    }
  }
}
```

jki (& kji):

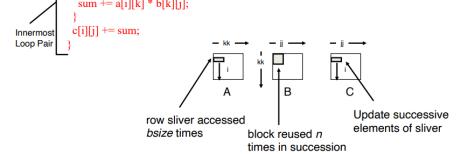
- 2 loads, 1 store
- misses/iter = 2.0

```
for (j=0; j<n; j++) {
  for (k=0; k<n; k++) {
    r = a[j][k];
    for (i=0; i<n; i++) {
      sum += r * b[k][i];
    }
  }
}
```

Blocked Matrix Multiply Analysis

- Innermost loop pair multiplies a $1 \times bsize$ sliver of A by a $bsize \times bsize$ block of B and accumulates into $1 \times bsize$ sliver of C
- Loop over i steps through n row slivers of A & C, using same B

```
for (i=0; i<n; i++) {
  for (j=jj; j< min(ij+bsize,n); j++) {
    sum = 0.0
    for (k=kk; k< min(kk+bsize,n); k++) {
      sum += a[i][k] * b[k][j];
    }
    c[i][j] += sum;
  }
}
```



Halide

Halide

- Halide is a domain specific language

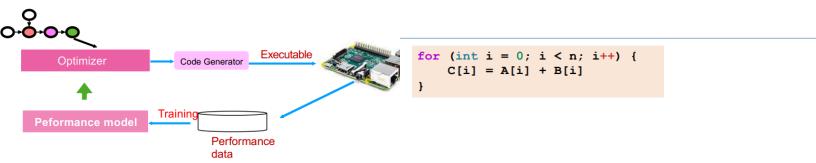
- It decoupled the optimized implementation

Program = algorithm + scheduling

- Enables programmers with different skills to collaborate
 - Algorithm experts and system experts
- Compilers are used to generate the parallel and optimized code

- Originally designed for image processing,
 - Now used by AI compilers for AI operator generation

Machine learning based optimizer



The machine learning model predict the performance within 1ms

Tensor-based abstraction

`C = matrix_multiply(A, B.T)`

```

m, n, h = tvm.var('m'), tvm.var('n'), tvm.var('h')
A = tvm.placeholder((m, h), name='A')
B = tvm.placeholder((n, h), name='B')

k = tvm.reduce_axis(0, h, name='k')
C = tvm.compute((m, n),
    lambda i, j: tvm.sum(A[i, k] * B[j, k], axis=k))
    
```

Inp
Output matrix shape
Computation rule:
$$C[i,j] = \sum_k A[i,k] \times B[j,k]$$

Scheduling example

```

C = tvm.compute((n, ), lambda i: A[i] + B[i])
s = tvm.create_schedule(C.op)
    
```

More tensor expressions

- Affine transformations

```

out = tvm.compute((n,m), lambda i, j: tvm.sum(data[i,k] * w[j,k], k))
out = tvm.compute((n,m), lambda i, j: out[i,j] + bias[i])
    
```

- Convolution

```

out = tvm.compute((c, h, w),
    lambda i, x, y: tvm.sum(data[kc,x+kx,y+ky] * w[i,kx,ky], [kx,ky]))
    
```

- ReLU

```

out = tvm.compute(shape, lambda *i: tvm.max(0, out(*i)))
    
```

Scheduling example

```

C = tvm.compute((n, ), lambda i: A[i] + B[i])
s = tvm.create_schedule(C.op)
    
```

```

xo, xi = s[C].split(s[C].axis[0], factor=32)
    
```

```

s[C].reorder(xi, xo)
    
```

```

s[C].bind(xo, tvm.thread_axis("blockIdx.x"))
    
```

```

s[C].bind(xi, tvm.thread_axis("threadIdx.x"))
    
```

Scheduling example

```

C = tvm.compute((n, ), lambda i: A[i] + B[i])
s = tvm.create_schedule(C.op)
    
```

```

xo, xi = s[C].split(s[C].axis[0], factor=32)
    
```

```

s[C].reorder(xi, xo)
    
```

```

s[C].bind(xo, tvm.thread_axis("blockIdx.x"))
    
```

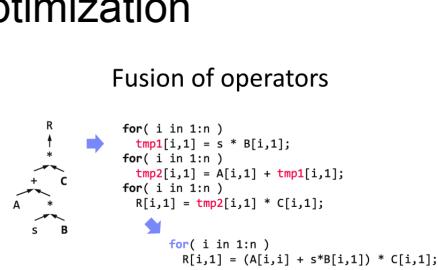
```

int i = threadIdx.x * 32 + blockIdx.x;
    
```

```

if (i < n) {
    C[i] = A[i] + B[i];
}
    
```

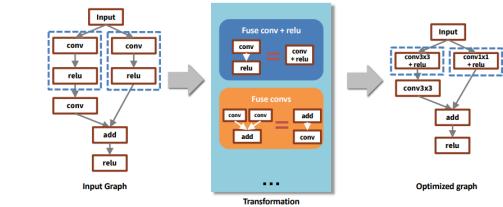
Transformation of Graphs



Fusion of operators

Graph Optimization

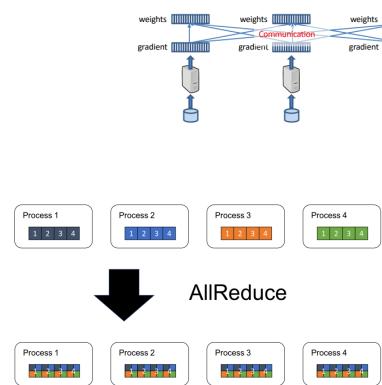
- Reduce the intermediate states which may cause memory references



Data Parallel Training

- In each iteration:

- each worker reads different mini-batch of size M, produces gradient
- communicate to produce aggregated gradient
- update weights

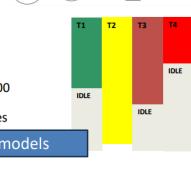
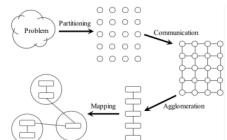


Distributed Training

Distributed training on 1000s cards

- Challenges
 - Communication reduction
 - Minimize the communications between between GPUs to allow high compute efficiency
 - Limited GPU memory
 - Intermediate state for backward computation
 - Load balance
 - Each card should have similar amount of work
 - Fault tolerance
 - MTBF is likely 10-100 hours, while training usually takes ~1000 hours
 - Quick recovery from system failures

10K – 100K cards for GPT4-like models

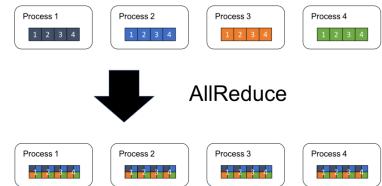
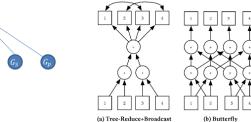


AllReduce to get global gradients

- We need to sum up all the gradients from different machines

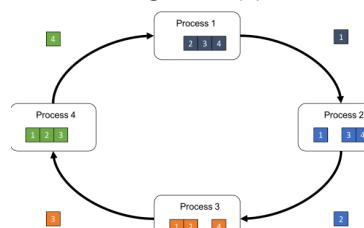
$$G = G_1 + G_2 + \dots + G_p$$

- Ways to do the summation:

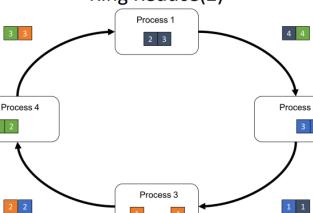


AllReduce

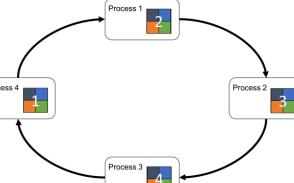
Ring Reduce(1)



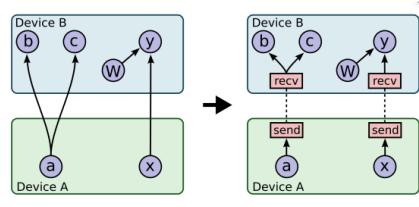
Ring Reduce(2)



Ring Reduce(3)

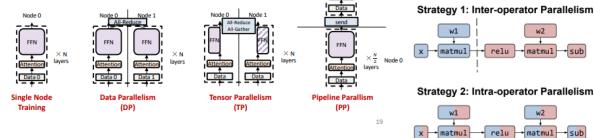


Model Parallel Training



Parallelisms for transformer

- Tensor model parallelism (split a layer)
- Pipeline parallelism (split between layers)

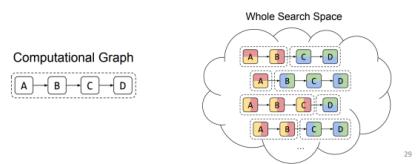


Challenges to computer systems :
How to partition the model and
choose parallelism?



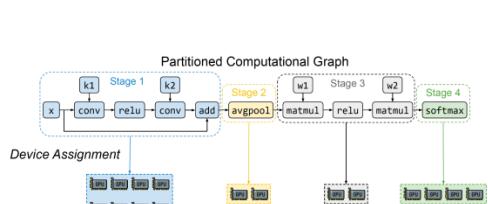
Can we find the good partition for a large model automatically

- Support both intra and inter-op parallelism, as well as pipeline parallelism
 - There are huge search space and it is a combinatorial optimization problem



Inter-op parallelism

Graph partition and device assignment within one compiler stage



Alpa

- Automatic parallelize the model training with just one line annotation

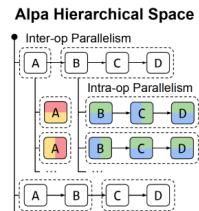
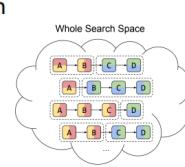
```

@alpa.parallelize
def train_step(model_state, batch):
    def loss_func(params):
        out = model_state.forward(params, batch["x"])
        return np.mean((out - batch["y"]) ** 2)

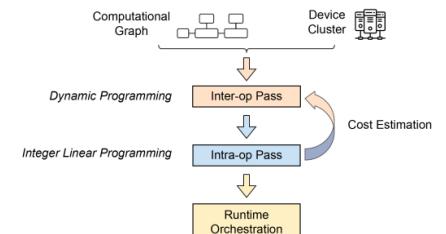
    grads = grad(loss_func)(state.params)
    new_model_state = model_state.apply_gradient(grads)
    return new_model_state

# A typical JAX training loop
model_state = create_train_state()
for batch in data_loader:
    model_state = train_step(model_state, batch)
    
```

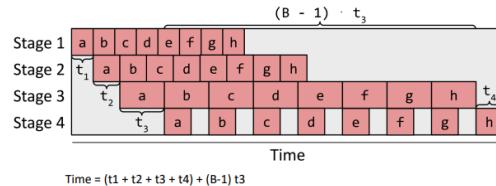
Two-level Hierarchical space to parallelism



Alpa Compiler phase



Inter-op parallelism



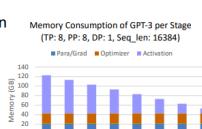
- Generalized objective function

$$T^* = \min_{\substack{s_1, \dots, s_S; \\ (n_1, m_1), \dots, (n_S, m_S)}} \left\{ \sum_{i=1}^S t_i + (B-1) \cdot \max_{1 \leq j \leq S} \{t_j\} \right\}$$

Memory Consumption

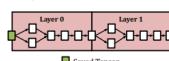
More on memory and load balance

- Large memory consumption
 - Optimizer
 - Parameters/Gradients
 - Activations
- Activations occupy huge memory
 - Recompute to save memory



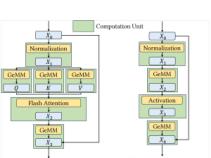
Recomputation

- Drop activations in the forward pass
- Recompute the activations in the backward pass
- Increased computation overhead



Computation Units

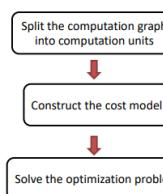
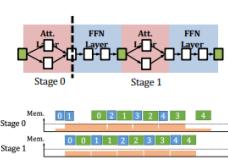
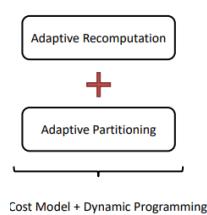
- Fine-grained partitioning
- Operators that require saving output tensors becomes the boundaries.
- Each Tensor bounds to one unit.



AdaPipe

AdaPipe

Adaptive Recomputation



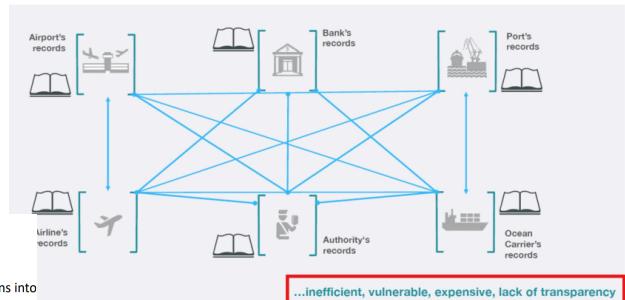
Requires interfaces between each pair of entities

Blockchain

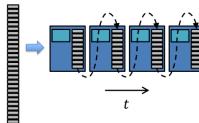
Motivation: Why correspondent banks are necessary— Lack of trust

Challenges of a global ledger

- Who can join this ledger?
 - Permissioned or permissionless
- Performance and availability
 - Global network latency and network partition
- Who can write what on this ledger?
 - How to avoid fake transactions?
 - How to avoid double spend?
 - How to avoid modification to previous transactions?
- Who should provide computers to maintain this ledger?
 - Why are they interested in doing this - Incentive



Performance and availability



- Batching transactions into blocks
 - Every 10 minutes 1MB block
- Large number of replicas
 - Every (full) node in the system has a copy

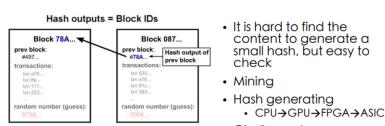
Who can join the ledger?

- Permissionless
 - Public chain
- Hyperledger
 - HYPERLEDGER
- Permissioned
 - Within a single organization
 - Private chain
 - Multiple party
 - Consortium Blockchain

Immutable data

- Crypto-hashing
 - Sha256
 - Single way function, sha256(byte array) = 256bits hash
 - Hard to construct byte array with a given hash
 - A block id which is the hash of the transactions in the block
 - A block contains the blockid of the previous block

Proof of work

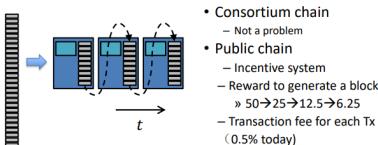


- It is hard to find the content to generate a small hash, but easy to check
- Mining
- Hash generating
 - CPU → GPU → FPGA → ASIC
- Challenge becomes more and more difficult

Choose transactions that:

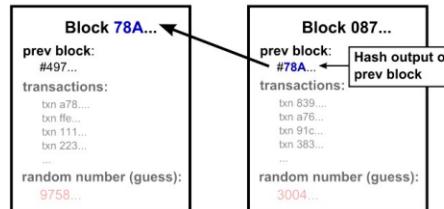
$\text{Sha256}(\text{prevID}, \text{transactions}, \text{guess}) < 00\ldots00010000$

(4) Why people are volunteering to maintain the blockchain systems

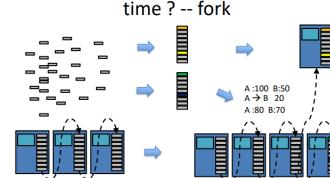


- Consortium chain
 - Not a problem
- Public chain
 - Incentive system
 - Reward to generate a block
 $\gg 50 \rightarrow 25 \rightarrow 12.5 \rightarrow 6.25$
 - Transaction fee for each Tx (0.5% today)

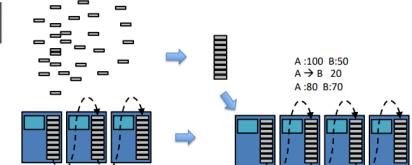
Hash outputs = Block IDs



What happens if two of the people solve the challenges at (almost) same time? -- fork



Consensus : Agree on the content of transactions in the block



What is a transaction?

What is a block?

What is a blockchain?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

What is a private key?

What is a digital signature?

What is a ledger?

What is a smart contract?

What is a consensus algorithm?

What is a proof of work?

What is a public key?

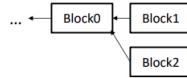
What is a private key?

What is a digital signature?

Forks in Nakamoto Consensus

Scalability:

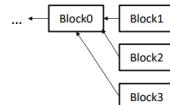
What if we run Nakamoto Consensus with **large/fast** blocks generation?



Large blocks → Longer propagation time
Fast block generation → Less time for propagation

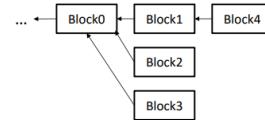
A new block is generated before the network propagates the previous block sufficiently

Forks in Nakamoto Consensus



More forks when the block generation is fast and

Forks in Nakamoto Consensus



Usually the system eventually converges on one fork by waiting for more blocks appended in the end

| Solution: | Projects: |
|------------|-----------------------------|
| Centralize | 7 nodes 21 nodes 1 node |

Wait! Why do we blockchain in the first place?

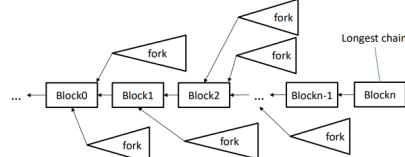
| Off-chain | Projects: |
|-----------|-------------------|
| | Lightning Network |

Great idea. But need large players to deposit on the chain.

| Sharding | Projects: |
|----------|-----------------|
| | Ethereum Casper |

Against CAP Theorem -- rely on economical mechanism to secure each shard.

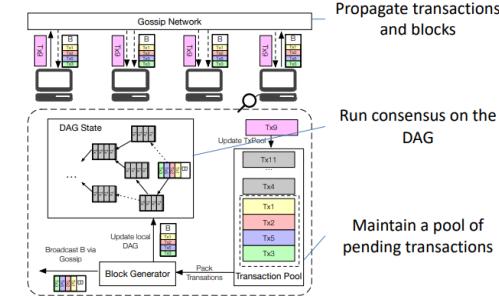
Forks in Nakamoto Consensus



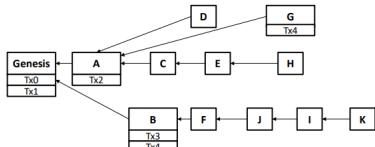
- Forks waste network/processing resources
- Downgrade safety
- Attacker needs less resource to beat the longest chain

Conflux

Conflux Architecture



Tx0: Mint 10 coins to X
Tx1: Mint 10 coins to Y
Tx2: X sends 8 to Y
Tx3: X sends 8 to Z
Tx4: Y sends 8 to Z

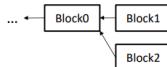


Each block has one outgoing parent edge to its parent block

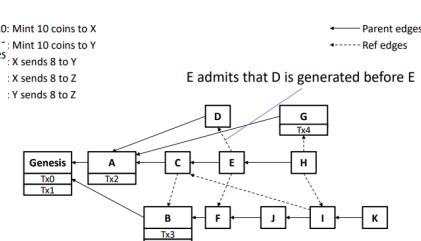
Parent edges form a tree

One Observation

- Bitcoin enforces a restrictive transaction total order at the generation time of each block:



- Blockchain transactions rarely conflict with each other and they can be serialized in any order
- Why not process non-conflicting transactions in all concurrent blocks?



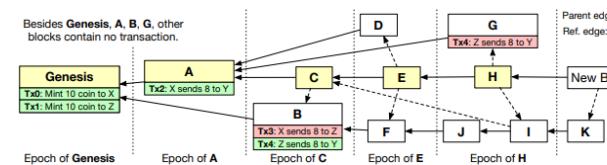
Each block may have multiple ref. edges

Ref. edges simply denote happens-before relationships

We still need a total order, how?

Conflux Consensus Ideas

- Optimistically process concurrent blocks without discarding any as forks
- Organize blocks into a **direct acyclic graphs (DAG)**
- Enable **fast/large block generation** without sacrificing security
 - First agree on a total order of all blocks
 - Assume transactions would not conflict each other
 - Then derives the transaction order from the agreed block order
 - Resolve transaction conflicts lazily



Computer Systems

Why is it difficult to build a system?

Scales of Computer Systems

- Programming / Data Structure
 - LoC (Lines of Code): From hundreds to millions
- Operating System / Architecture
 - Cores: from one to hundreds
- Network
 - Nodes: from two to millions
- Web Service
 - Clients: from tens to millions

Complexity of Computer Systems

- Hard to define; symptoms:
 - Large number of components
 - Large number of connections
 - Irregular
 - No short description
 - Many people required to design/maintain
- Technology rarely the limit
 - Indeed tech opportunity is the problem
 - Limit is usually designers' understanding

Common Problems of Systems

Problem Types

- **Emergent properties** (surprise!)
 - The properties that are not considered at design time
- **Propagation of effects**
 - Small change -> big effect
- **Incommensurate scaling**
 - Design for small model may not scale
- **Trade-offs**
 - Waterbed effect

Waterbed Effect

- Pushing down on a problem at one point
- Causes another problem to pop up somewhere else

Coping with Complexity

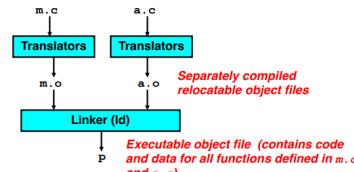
M.A.L.H

- **Modularity**
 - Split up system
 - Consider separately
- **Layering**
 - Gradually build up capabilities
- **Abstraction**
 - Interface/Hiding
 - Avoid propagation of effects
- **Hierarchy**
 - Reduce connections
 - Divide-and-conquer

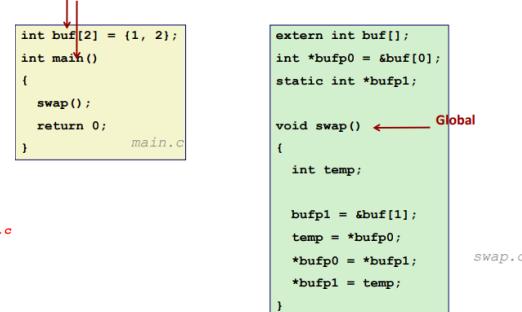
Layering

- Goal
 - Reduce module interconnections even more
- How to do it
 - Build a set of mechanisms first(a lower layer)
 - Use them to create a different complete set of mechanisms (an upper layer)
- General rule: A module in one layer only interacts with:
 - Its peers in the same layer, and
 - Modules in the next lower layer / next higher layer

Use linker to support modularity in software design



How pointers break modular boundaries



Unbounded Composition

- Two properties of computer systems
 - 1. Mostly digital
 - 2. Controlled by software
- Both relax the limits on complexity arising from physical laws in other systems

Prep. Exam

1. (15 points) In a storm application, a Spout with <sid1> sends a tuple to Bolt1 whose stid is 0010. It's acker is Ack1. After the execution of Bolt1, it sends two tuples to Bolt2 and Bolt3 respectively.

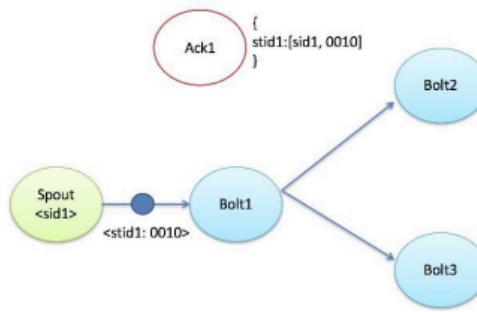
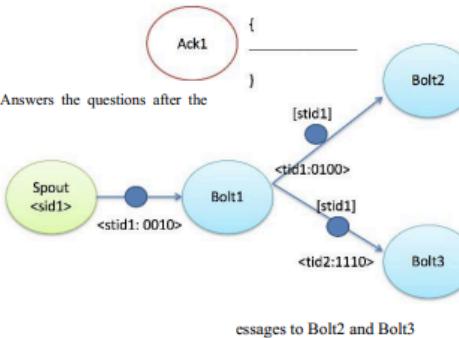


Figure 1.1 Spout sid1 sends a tuple to Bolt 1

3. (15points) For the following MPI algorithm. Answers the questions after the source code.



messages to Bolt2 and Bolt3

```

#include <stdlib.h>
#include <stdio.h>
#include <math.h>
#include <time.h>
#include "mpi.h"
int main ( int argc, char *argv[] );
int prime_number ( int n, int id, int p );

int main ( int argc, char *argv[] )
{
    int i, id, ierr, n, n_factor, n_hi, n_lo, p, primes, primes_part;
    double wtime;
    n_lo = 1;
    n_hi = 262144;
    n_factor = 2;
    ierr = MPI_Init ( &argc, &argv );
    ierr = MPI_Comm_size ( MPI_COMM_WORLD, &p );
    ierr = MPI_Comm_rank ( MPI_COMM_WORLD, &id );
    ierr = MPI_Wtime ( &wtime );
    primes_part = prime_number ( n, id, p );
    ierr = MPI_Reduce ( &primes_part, &primes, 1, MPI_INT,
    MPI_SUM, 0, MPI_COMM_WORLD );
    if ( id == 0 ) {
        wtime = MPI_Wtime ( ) - wtime;
        printf ( "%8d %8d %14fn", n, primes, wtime );
    }
    ierr = MPI_Finalize ( );
    return 0;
}

int prime_number ( int n, int id, int p )
{
    Purpose:
        PRIME_NUMBER returns the number of primes between 1 and N.
    Parameters:
        Input, int N, the maximum number to check.
        Input, int ID, the ID of this process,
        between 0 and P-1.
        Input, int P, the number of processes.
        Output, int PRIME_NUMBER, the number of prime numbers up to N.
    {
        int i, total, prime;
        total = 0;

        for ( i = 2 + id; i <= n; i = i + p )
            prime = 1;
            for ( j = 2; j < i; j++ )
                if ( ( i % j ) == 0 )
                    prime = 0;
                    break;
                }
            total = total + prime;
            }
        return total;
}
  
```

2. (15 points) In graph computing frameworks, sometimes it is useful to sort vertex with the degree, such as the bottom-up algorithm for BFS. Now we need to write a spark program for this task.

Input:

An unsorted edge list of a directed graph with format(source_vertex_id, destination_vertex_id), the vertex id is a unsigned integer less than 4 billion. The edge list is in text format, one edge per line. For example.

1,2
2,4
3,5
1,4

Output:

Vertex id sorted by in-degree with descending order. If two vertices have the same degree, they are sorted by vertex id with ascending order. The output of the above example is:

4
2
5
1
3

Question:

Write a spark program to perform the task.

```

lines = sc.textFile(input)
edges = lines.map(lambda line: tuple(line.split(',')))
degree_v = edges.map(lambda e: (e[1], 1)).reduce(add).sortBy(lambda x: (-x[1],
x[0])).collect()
for _d, v in degree_v:
    print(v)
  
```

- a) Please fill the blanks in the Figure 1.2 to show the record of stid1. The format example is shown in Figure 1.1.

After Bolt1 send ack to Ack1, record will be {stid1: [stid1: 1010]}

- b) How would the system detect that a message is not processed completely and need to be reprocessed?

If a bolt crashes, or emits a failure not ack, or the tuple reaches timeout.

- c) What is the feature of the Spout to enable the feature of message reprocessing? What is the typical software used to support this feature?

Replayable source. Queue (Kafka).

Question 1: Fill in the blanks in the first page.

Question 2: Is this a load balanced algorithm? How many numbers are processed by each process with number n and process number p(roughly) in one iteration?

Yes. (n-2)/p.

Question3: Any ideas to optimize this code?

Many ways. For example, for (j = 2; j < i; j++) -> for (j = 2; j*j < i; j++) .

